*SIGMOD 2015, Melbourne, Victoria, Australia*

# Rethinking SIMD Vectorization for In-Memory Databases

**Orestis Polychroniou**
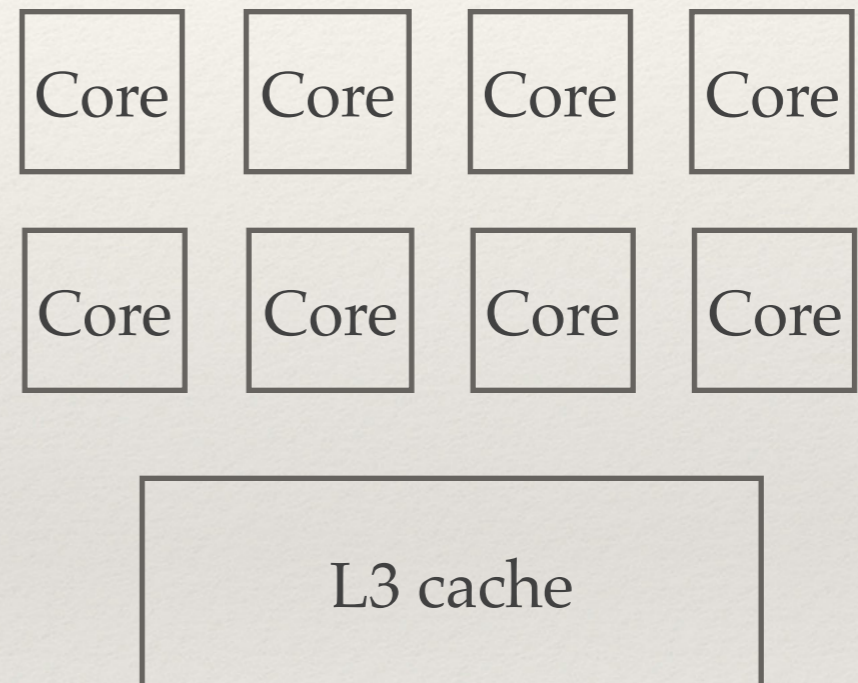*Columbia University*

**Arun Raghavan**
*Oracle Labs*

**Kenneth A. Ross**
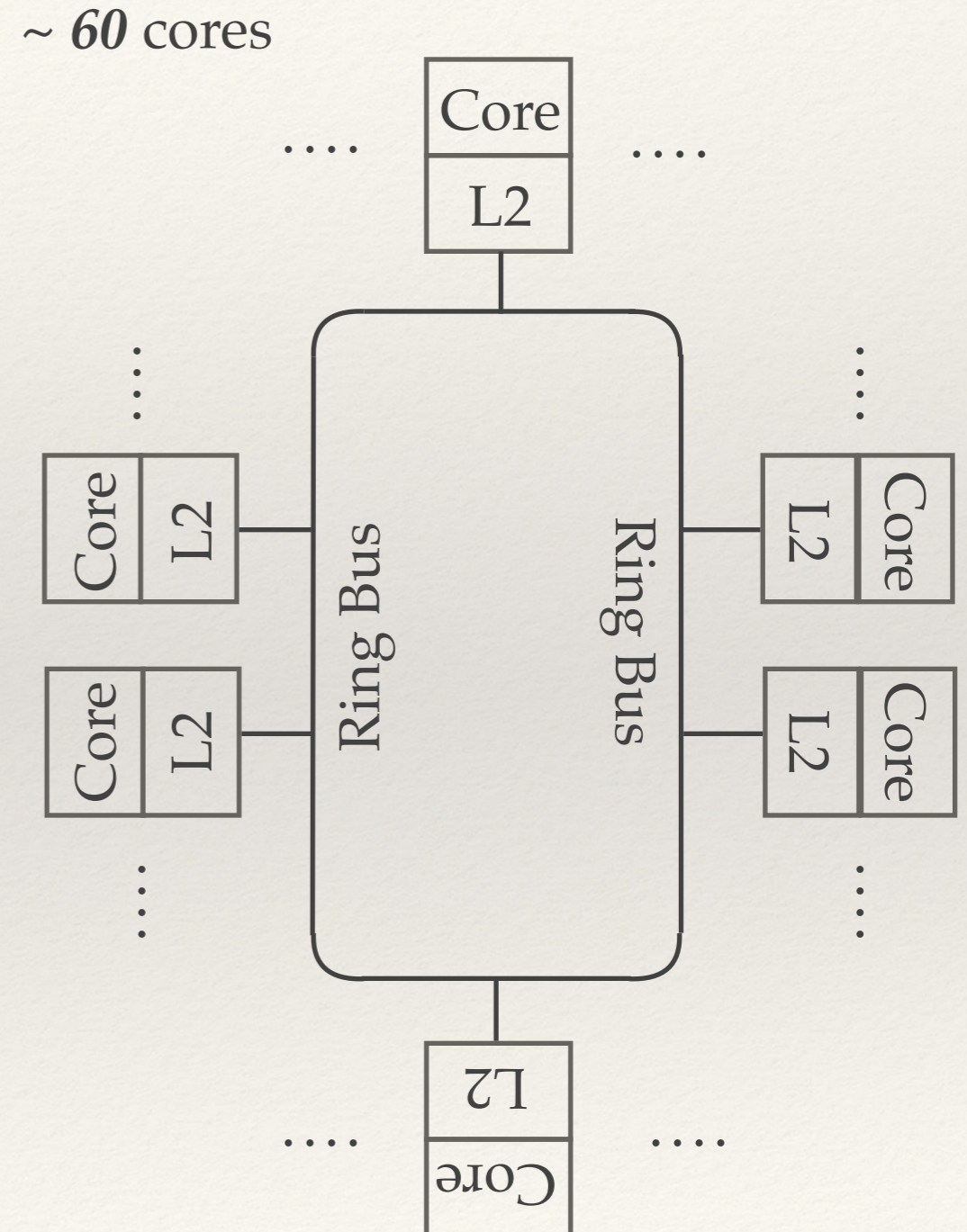*Columbia University*

# Latest Hardware Designs

❖ Mainstream multi-core CPUs

    ❖ Use *complex* cores  (e.g. Intel Haswell)

       ❖ Massively *superscalar*

       ❖ Aggressively *out-of-order*

    ❖ Pack *few* cores per chip

       ❖ High *power & area* per core

*2 – 18* cores

| Core | Core | Core | Core |
|------|------|------|------|
| Core | Core | Core | Core |

L3 cache

# Latest Hardware Designs

- Mainstream multi-core CPUs
  - Use *complex* cores  (e.g. Intel Haswell)
    - Massively *superscalar*
    - Aggressively *out-of-order*
  - Pack *few* cores per chip
    - High *power & area* per core

- Many Integrated Cores  (MIC)
  - Use *simple* cores  (e.g. Intel P54C)
    - *In-order* & *non-superscalar*  (for SIMD)
  - *Augment* SIMD to bridge the gap
    - Increase SIMD *register size*
    - More *advanced* SIMD instructions
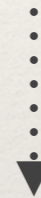  - Pack *many* cores per chip
    - Low *power & area* per core

*~ 60* cores

# SIMD & Databases

- ❖ Automatic vectorization
  - ❖ Works for *simple* loops only
    - ❖ *Rare* in database operators

*what is **SIMD** ?*

$$a + b = c$$

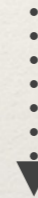| $a1$ | $a2$ | $a3$ | $a4$ | $\cdots$ | $aW$ |
| --- | --- | --- | --- | --- | --- |
| $+$ | $+$ | $+$ | $+$ | | $+$ |
| $b1$ | $b2$ | $b3$ | $b4$ | $\cdots$ | $bW$ |
| $=$ | $=$ | $=$ | $=$ | | $=$ |
| $c1$ | $c2$ | $c3$ | $c4$ | $\cdots$ | $cW$ |

# SIMD & Databases

- ❖ Automatic vectorization
  - ❖ Works for *simple* loops only
    - ❖ *Rare* in database operators

- ❖ Manual vectorization
  - ❖ *Linear* access operators
    - ❖ Predicate evaluation
    - ❖ Compression
  - ❖ *Ad-hoc* vectorization
    - ❖ Sorting (e.g. merge-sort, comb-sort, …)
    - ❖ Merging (e.g. 2-way, bitonic, …)
  - ❖ Generic vectorization
    - ❖ *Multi-way* trees
    - ❖ *Bucketized* hash tables

*what is SIMD ?*

$$a + b = c$$

| $a1$ | $a2$ | $a3$ | $a4$ | $\cdots$ | $aW$ |
|------|------|------|------|----------|------|
| $+$  | $+$  | $+$  | $+$  |          | $+$  |
| $b1$ | $b2$ | $b3$ | $b4$ | $\cdots$ | $bW$ |
| $=$  | $=$  | $=$  | $=$  |          | $=$  |
| $c1$ | $c2$ | $c3$ | $c4$ | $\cdots$ | $cW$ |

# Contributions & Outline

❖ *Full* vectorization

    ❖ From $O(f(n))$ scalar to $O(f(n)/W)$ vector operations

        ❖ *Random* accesses excluded

    ❖ *Principles* for good (*efficient*) vectorization

        ❖ Reuse fundamental *operations* across multiple vectorizations

        ❖ Favor *vertical vectorization* by processing different input data per lane

        ❖ Maximize lane *utilization* by executing different things per lane subset

# Contributions & Outline

- *Full* vectorization

  - From $O(f(n))$ scalar to $O(f(n)/W)$ vector operations

    - *Random* accesses excluded

  - *Principles* for good (*efficient*) vectorization

    - Reuse fundamental *operations* across multiple vectorizations

    - Favor *vertical vectorization* by processing different input data per lane

    - Maximize lane *utilization* by executing different things per lane subset

- Vectorize *basic* operators

  - Selection scans

  - Hash tables

  - Partitioning

# Contributions & Outline

❖ *Full* vectorization

   ❖ From $O(f(n))$ scalar to $O(f(n)/W)$ vector operations

      ❖ *Random* accesses excluded

   ❖ *Principles* for good (*efficient*) vectorization

      ❖ Reuse fundamental *operations* across multiple vectorizations

      ❖ Favor *vertical vectorization* by processing different input data per lane

      ❖ Maximize lane *utilization* by executing different things per lane subset

❖ Vectorize *basic* operators & build *advanced* operators  (in memory)

   ❖ Selection scans

   ❖ Hash tables
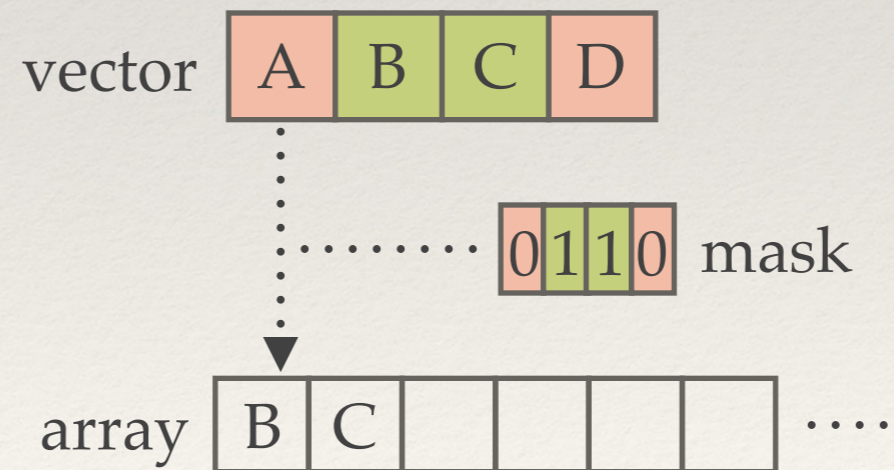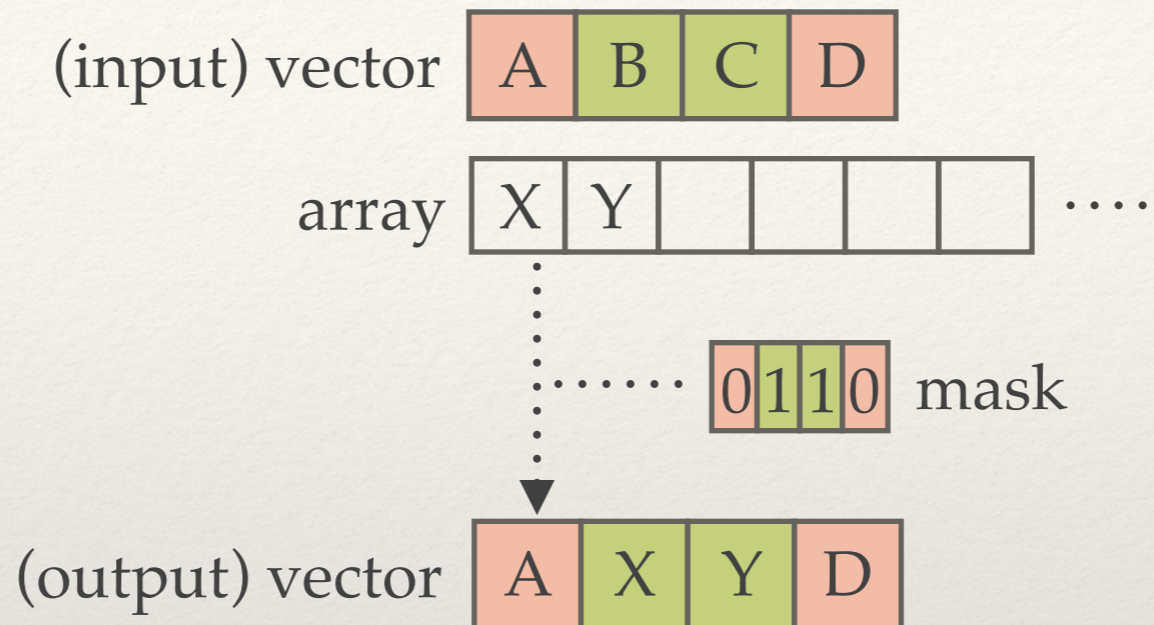
   ❖ Partitioning

   ❖ Sorting

   ❖ Joins

# Contributions & Outline

❖ *Full* vectorization

  ❖ From $O(f(n))$ scalar to $O(f(n)/W)$ vector operations

    ❖ *Random* accesses excluded

  ❖ *Principles* for good (*efficient*) vectorization

    ❖ Reuse fundamental *operations* across multiple vectorizations

    ❖ Favor *vertical vectorization* by processing different input data per lane

    ❖ Maximize lane *utilization* by executing different things per lane subset

❖ Vectorize *basic* operators & build *advanced* operators (in memory)

  ❖ Selection scans

  ❖ Hash tables

  ❖ Partitioning

  ❖ Sorting

  ❖ Joins

❖ Show impact of *good* vectorization

  ❖ On *software* (database system) & *hardware* design
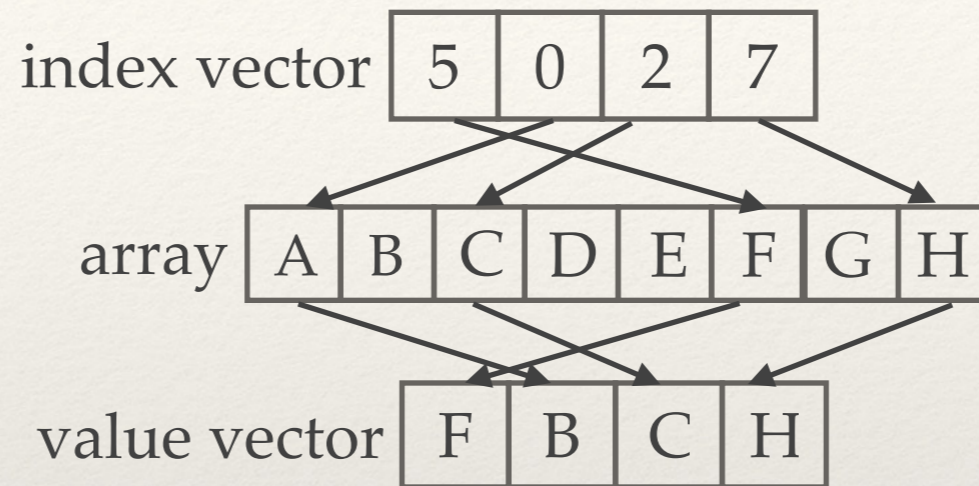
# Fundamental Vector Operations

- Selective load

- Selective store

(input) vector  | A | B | C | D |
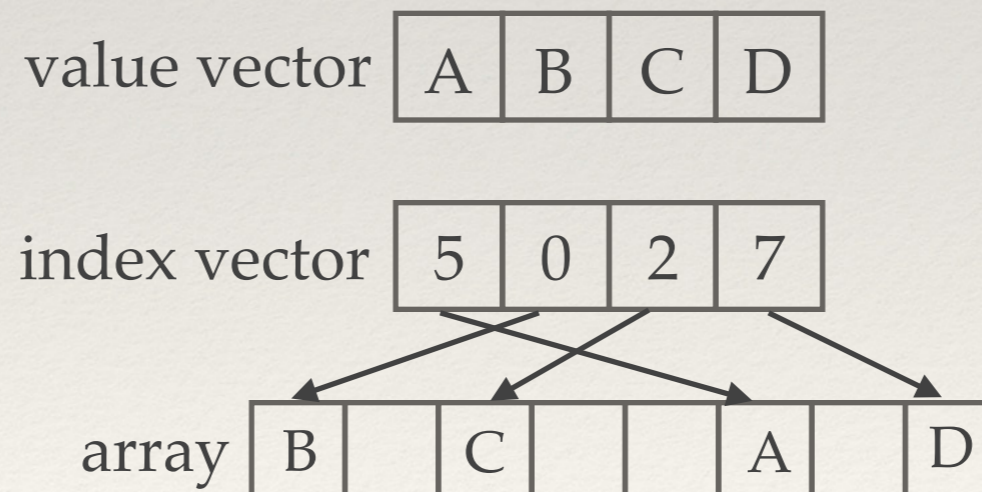
array | X | Y | | | | | ....

0 1 1 0  mask

(output) vector | A | X | Y | D |

vector | A | B | C | D |

0 1 1 0  mask

array | B | C | | | | | ....

# Fundamental Vector Operations

- ❖ Selective load

- ❖ Selective store

- ❖ (Selective) gather

- ❖ (Selective) scatter

index vector | 5 | 0 | 2 | 7

array | A | B | C | D | E | F | G | H

value vector | F | B | C | H

value vector | A | B | C | D

index vector | 5 | 0 | 2 | 7

array | B | | C | | A | | D

# Selection Scans

- ❖ Scalar
  - ❖ Branching
    - ❖ *Branch* to store qualifiers
    - ❖ Affected by *branch misses*
  - ❖ Branchless
    - ❖ *Eliminate* branching
    - ❖ Use conditional *flags*

# Selection Scans

- Scalar
  - Branching
    - *Branch* to store qualifiers
    - Affected by *branch misses*
  - Branchless
    - *Eliminate* branching
    - Use conditional *flags*

- Vectorized
  - Simple
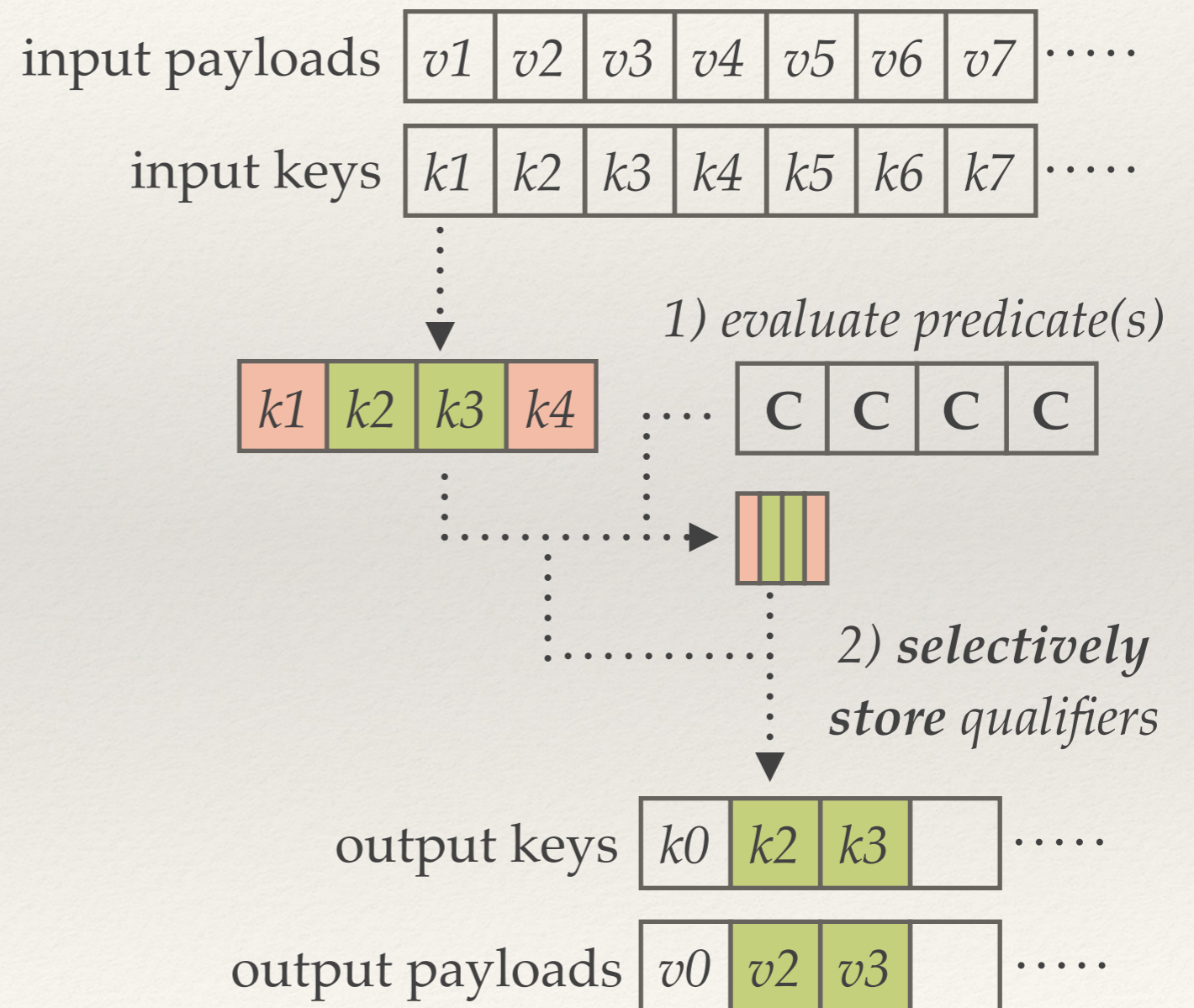    - Evaluate predicates (in SIMD)
    - *Selectively store* qualifiers
    - "Early" materialized
  - Advanced
    - "Late" materialized (in the paper …)

*select … where column > **C** …*

input payloads | v1 | v2 | v3 | v4 | v5 | v6 | v7 | · · · · ·

input keys | k1 | k2 | k3 | k4 | k5 | k6 | k7 | · · · · ·

*1) evaluate predicate(s)*

| k1 | k2 | k3 | k4 | · · · · | **C** | **C** | **C** | **C** |

*2) **selectively** store qualifiers*

output keys | k0 | k2 | k3 | | · · · · ·

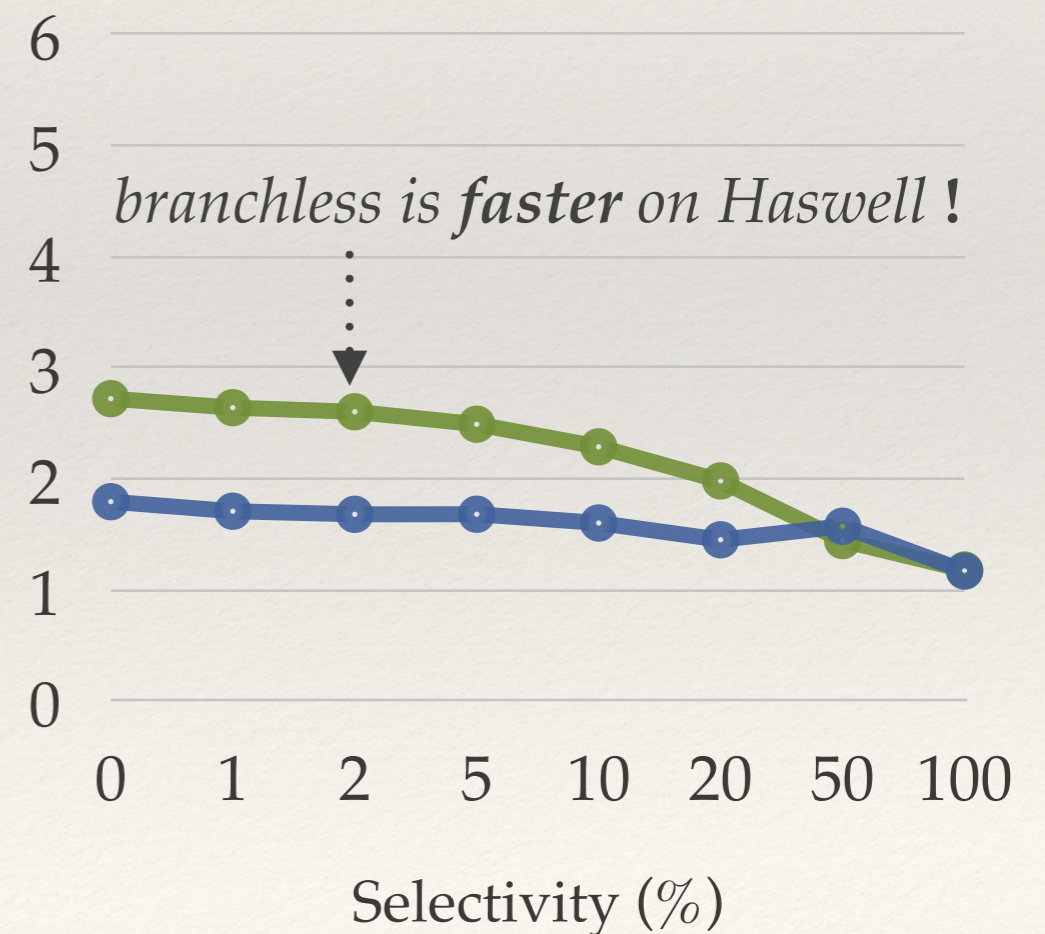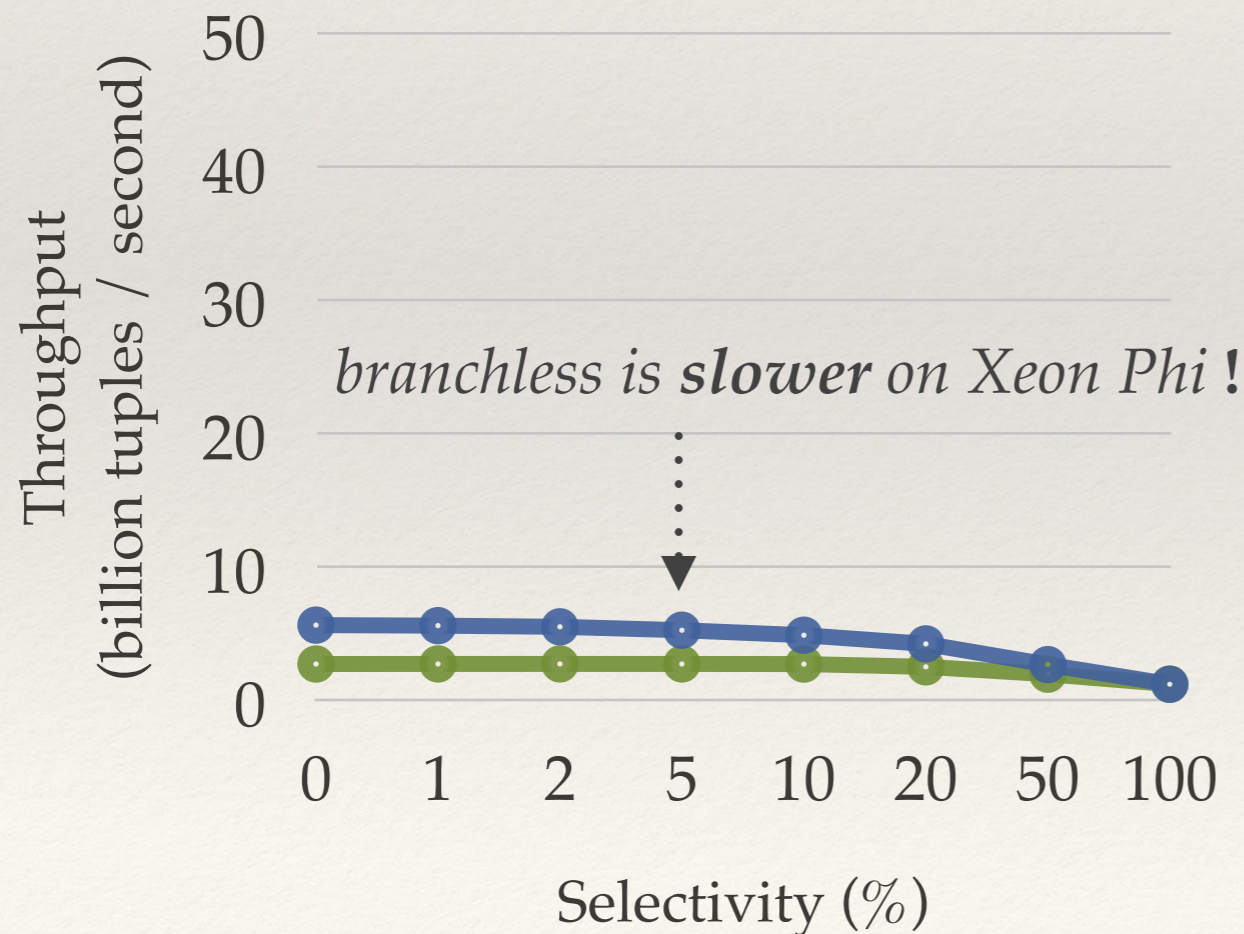output payloads | v0 | v2 | v3 | | · · · · ·

# Selection Scans

Xeon Phi 7120P (**MIC**)

○ Scalar (branching)
○ Scalar (branchless)

Xeon E3-1275v3 (**Haswell**)
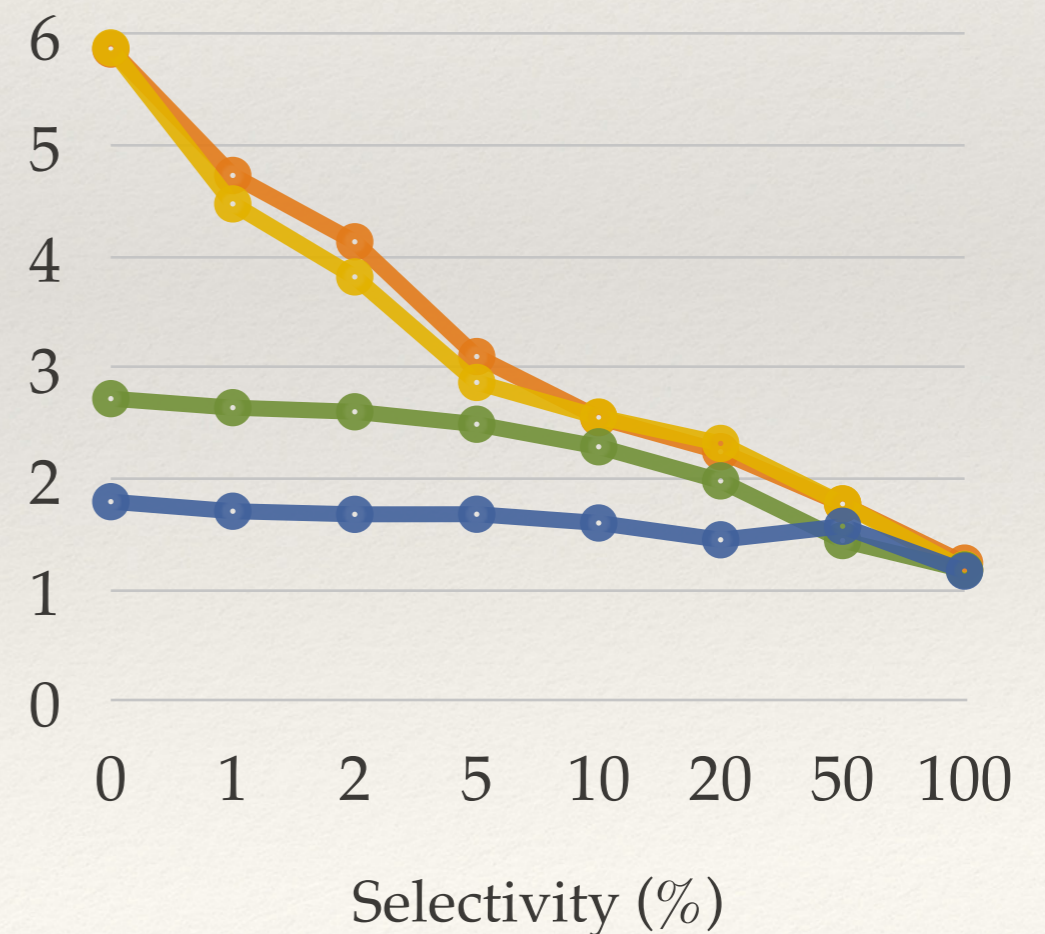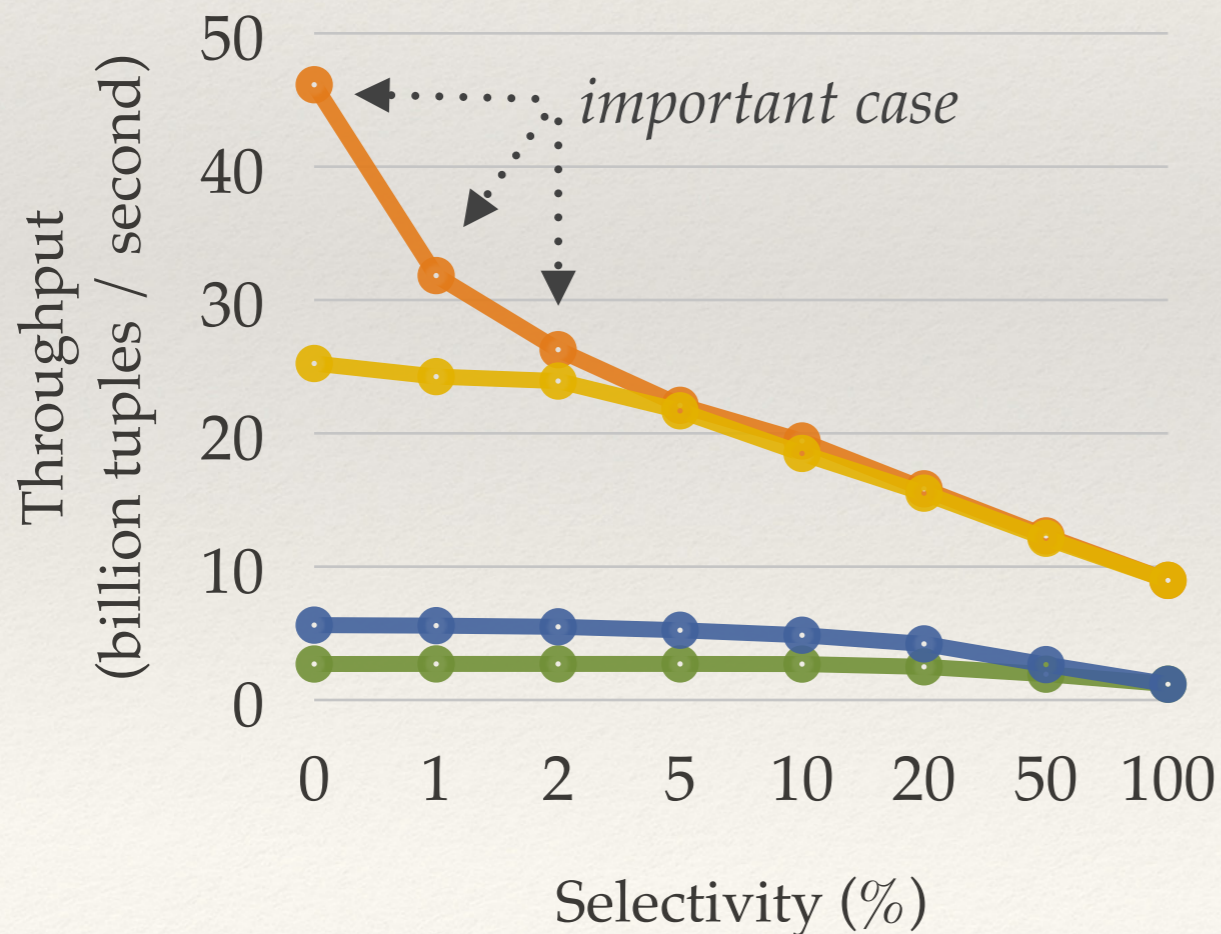
○ Scalar (branching)
○ Scalar (branchless)

*branchless is **slower** on Xeon Phi !*

*branchless is **faster** on Haswell !*

Throughput (billion tuples / second)

Selectivity (%)

Selectivity (%)

# Selection Scans

# Hash Tables

- Scalar hash tables
  - Branching or branchless
    - 1 *input* key at a time
    - 1 *table* key at a time

linear probing
hash table

input
key

hash
index

$k$ ┈┈▶ $h$ ┈┈┈▶

keys

payloads

*probe > 1
buckets ?*

# Hash Tables

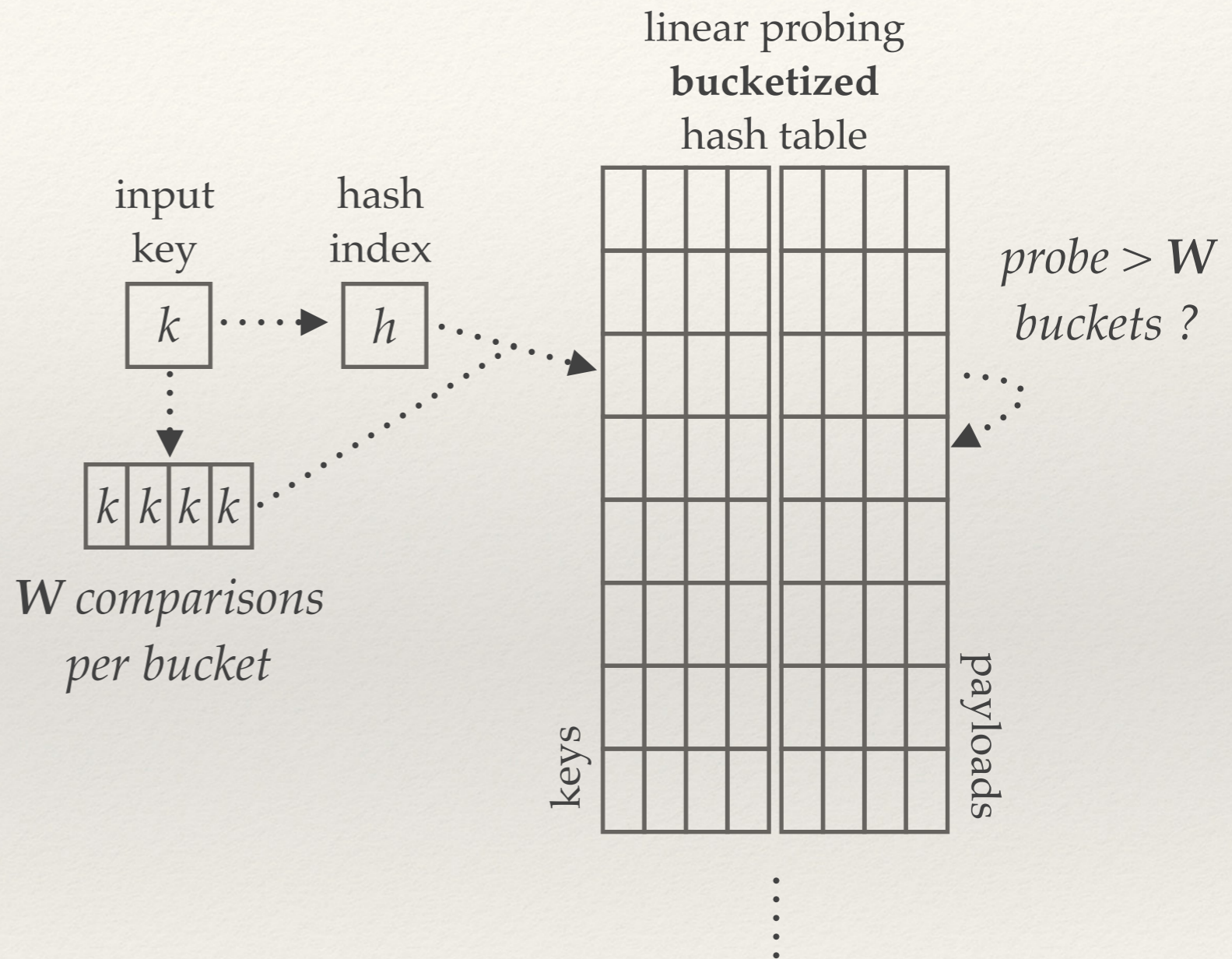- ❖ Scalar hash tables
  - ❖ Branching or branchless
    - ❖ 1 *input* key at a time
    - ❖ 1 *table* key at a time

- ❖ Vectorized hash tables
  - ❖ *Horizontal* vectorization
    - ❖ Proposed on *previous work*
    - ❖ 1 *input* key at a time
    - ❖ W *table* keys per input key
    - ❖ *Load* bucket with *W* keys

linear probing
**bucketized**
hash table

input
key

hash
index

$k$ ⋯▶ $h$

$k$ $k$ $k$ $k$

*W comparisons*
*per bucket*

keys

payloads

*probe > **W***
*buckets ?*

# Hash Tables

- Scalar hash tables

  - Branching or branchless

    - 1 *input* key at a time

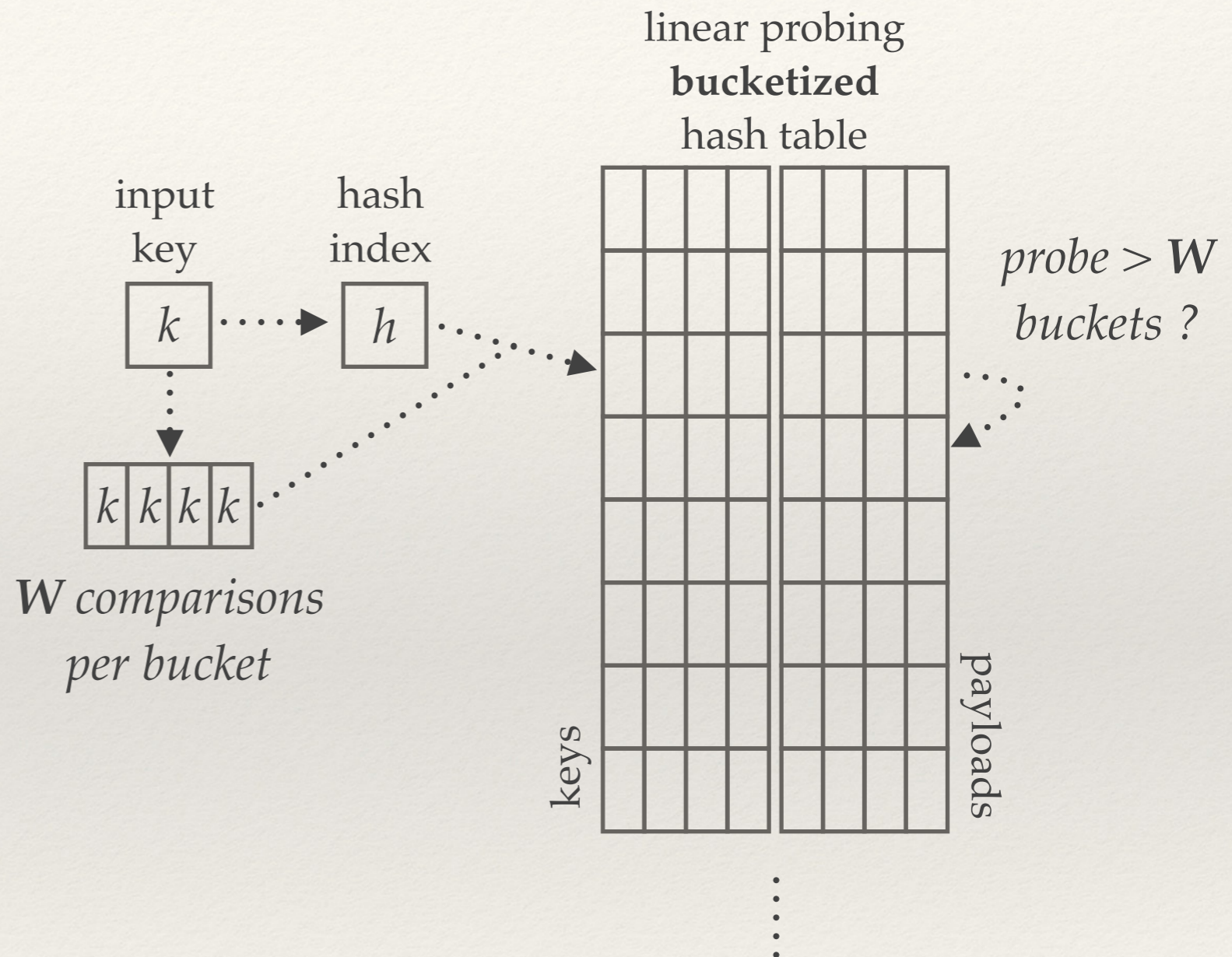    - 1 *table* key at a time

- Vectorized hash tables

  - *Horizontal* vectorization

    - Proposed on *previous work*

    - 1 *input* key at a time

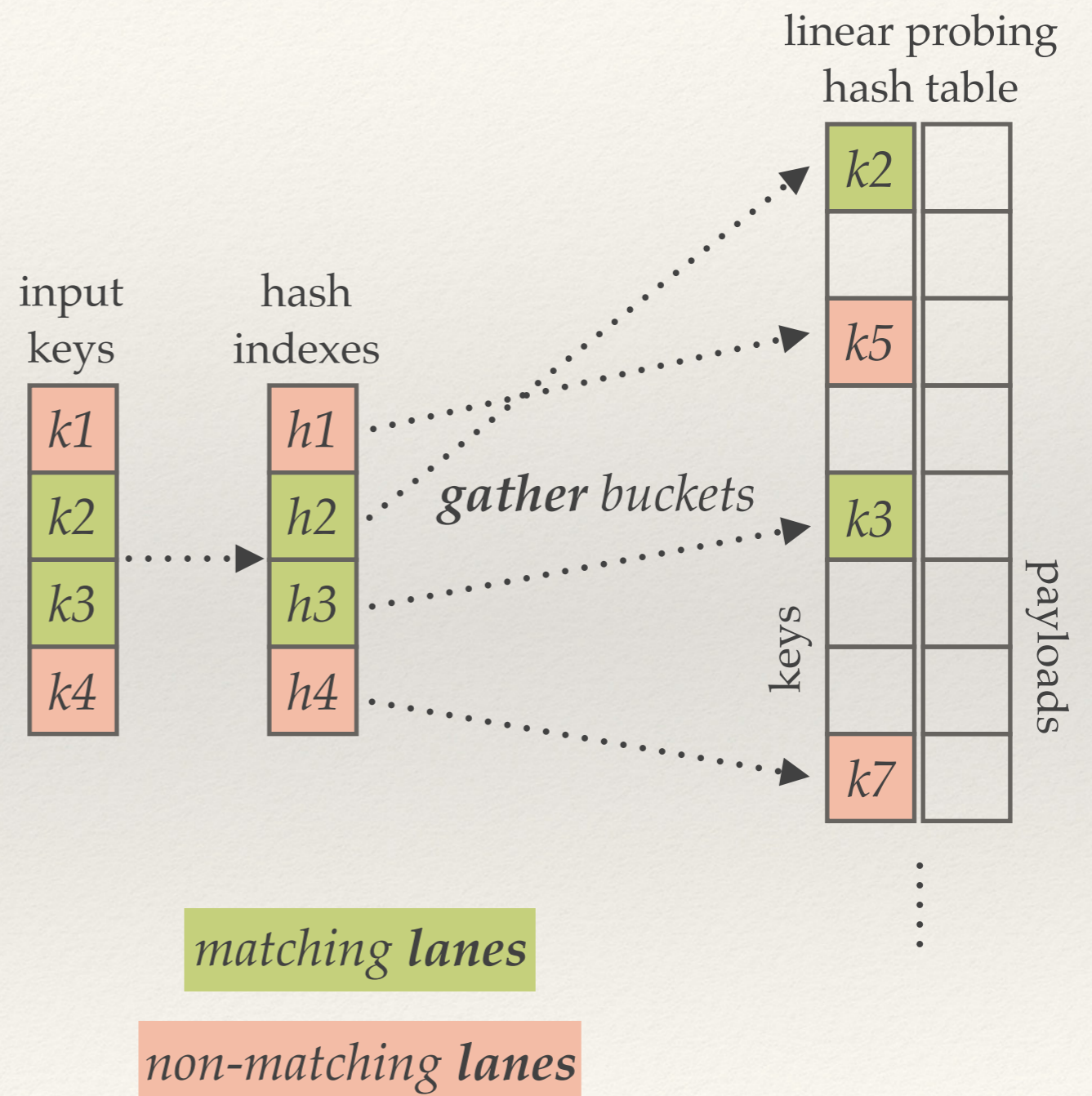    - W *table* keys per input key

    - *Load* bucket with *W* keys

  - However …

    - *W* are too many comparisons

    - No advantage of larger SIMD

linear probing
**bucketized**
hash table

input
key

hash
index

$k$

$h$

$k\ k\ k\ k$

*W comparisons
per bucket*

*probe > **W**
buckets ?*

keys

payloads

# Hash Tables

- Vectorized hash tables
  - *Vertical* vectorization
    - W *input* keys at a time
    - 1 *table* keys per input key
    - *Gather* buckets
  - Probing *(linear probing)*
    - *Store* matching lanes



input keys

hash indexes

linear probing hash table

*gather* buckets

keys

payloads

k1 h1 k2

k2 h2 k5

k3 h3 k3

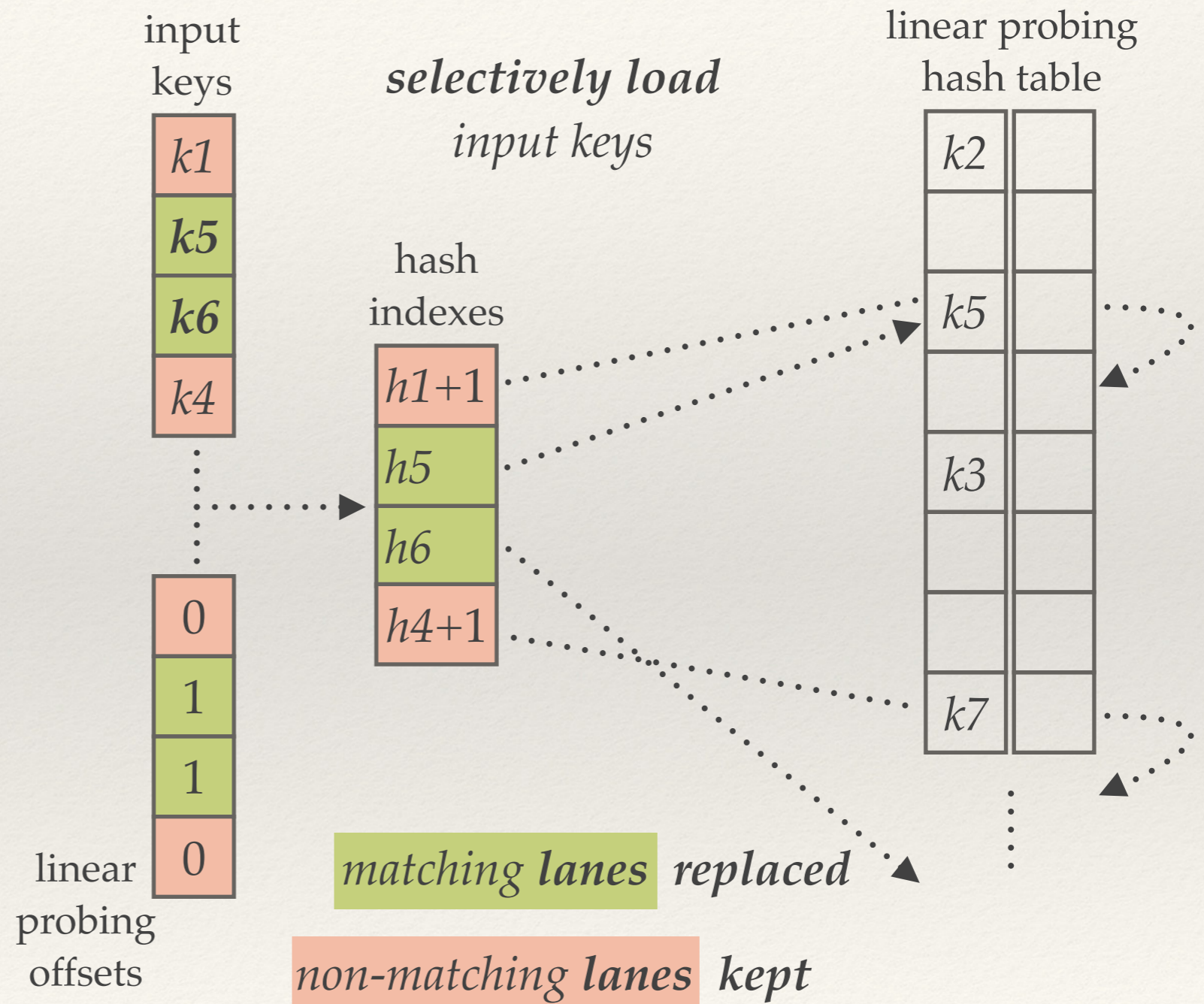k4 h4 k7

matching **lanes**

non-matching **lanes**

# Hash Tables

- Vectorized hash tables
  - *Vertical* vectorization
    - W *input* keys at a time
    - 1 *table* keys per input key
    - *Gather* buckets
  - Probing *(linear probing)*
    - *Store* matching lanes
    - *Replace* finished lanes
    - *Keep* unfinished lanes



input keys

| k1 |
| k5 |
| k6 |
| k4 |

*selectively load*
input keys

hash indexes

| h1+1 |
| h5 |
| h6 |
| h4+1 |

linear probing hash table

| k2 |   |
|    |   |
| k5 |   |
|    |   |
| k3 |   |
|    |   |
|    |   |
| k7 |   |

linear probing offsets

| 0 |
| 1 |
| 1 |
| 0 |

*matching lanes replaced*

*non-matching lanes kept*

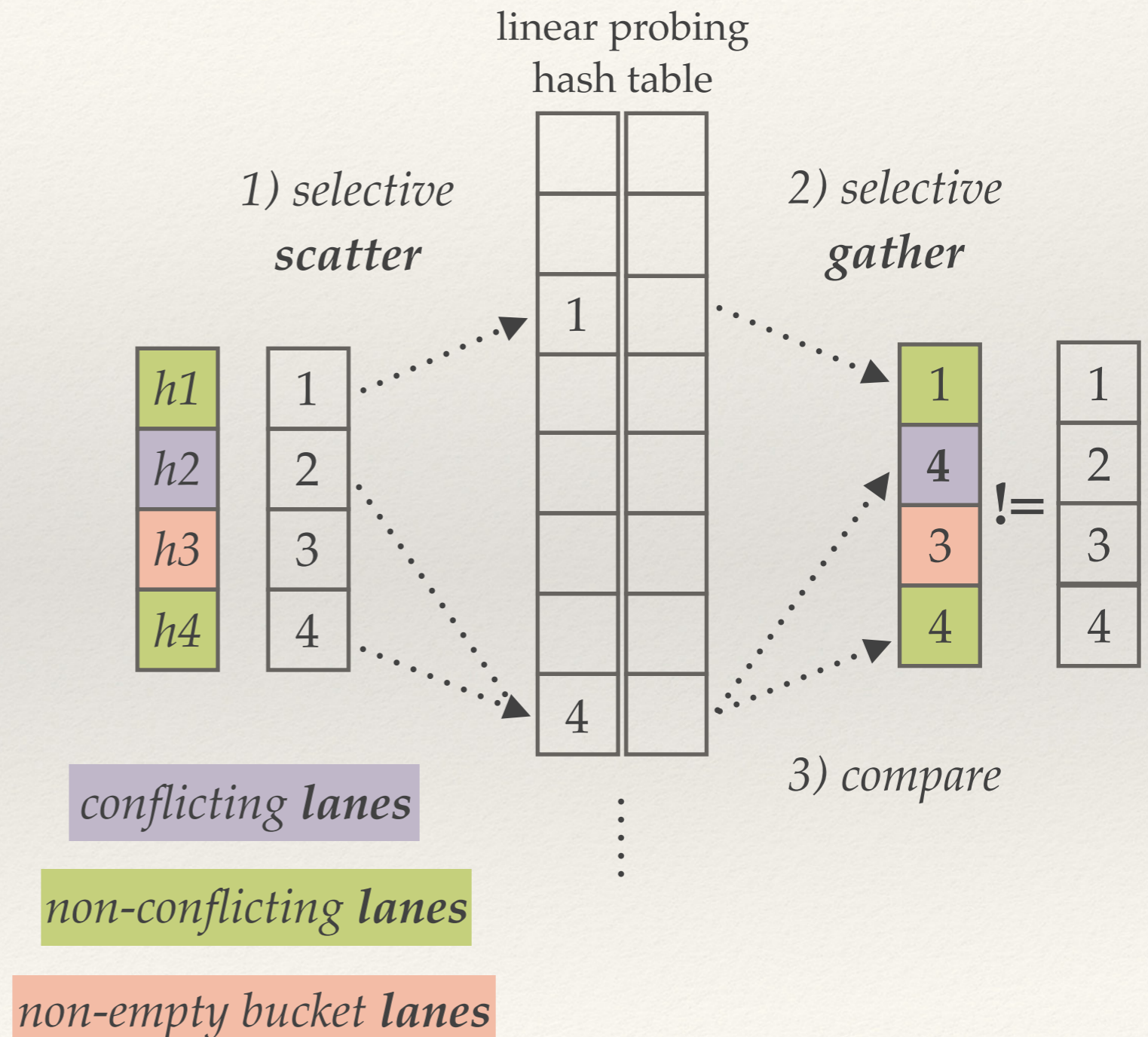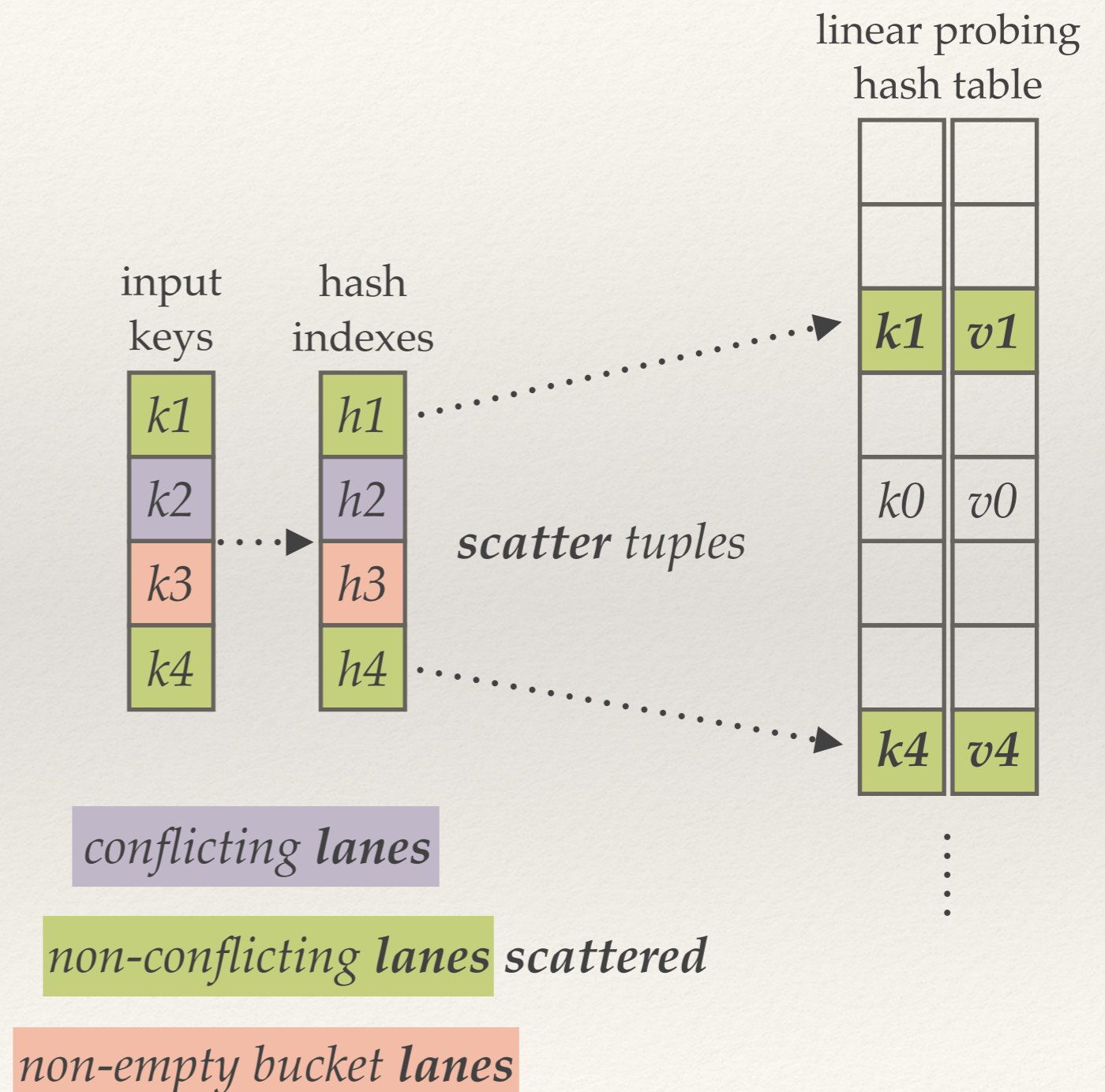# Hash Tables

❖ Vectorized hash tables

   ❖ *Vertical* vectorization

      ❖ W *input* keys at a time

      ❖ 1 *table* keys per input key

      ❖ *Gather* buckets

   ❖ Probing  *(linear probing)*

      ❖ *Store* matching lanes

      ❖ *Replace* finished lanes

      ❖ *Keep* unfinished lanes

   ❖ Building  *(linear probing)*

      ❖ *Keep* non-empty lanes

linear probing
hash table

input
keys

hash
indexes

| k1 | h1 |
| k2 | h2 |
| k3 | h3 |
| k4 | h4 |

*gather* buckets

null

k0

null

*empty bucket lanes*
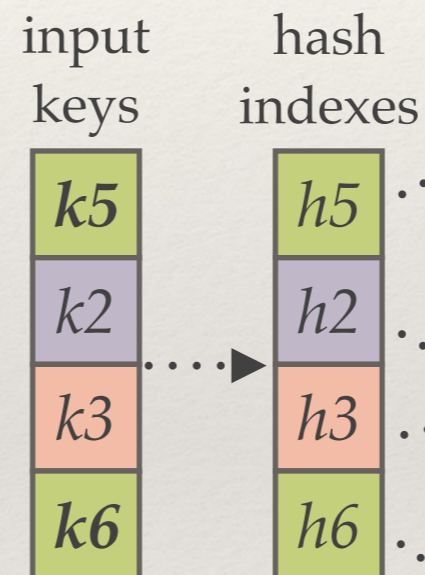
*non-empty bucket lanes*

# Hash Tables

- Vectorized hash tables

  - *Vertical* vectorization

    - W *input* keys at a time
    - 1 *table* keys per input key
    - *Gather* buckets

  - Probing *(linear probing)*

    - *Store* matching lanes
    - *Replace* finished lanes
    - *Keep* unfinished lanes

  - Building *(linear probing)*

    - *Keep* non-empty lanes
    - *Detect* scatter *conflicts*



linear probing hash table

1) selective **scatter**

2) selective **gather**

3) compare

conflicting **lanes**

non-conflicting **lanes**
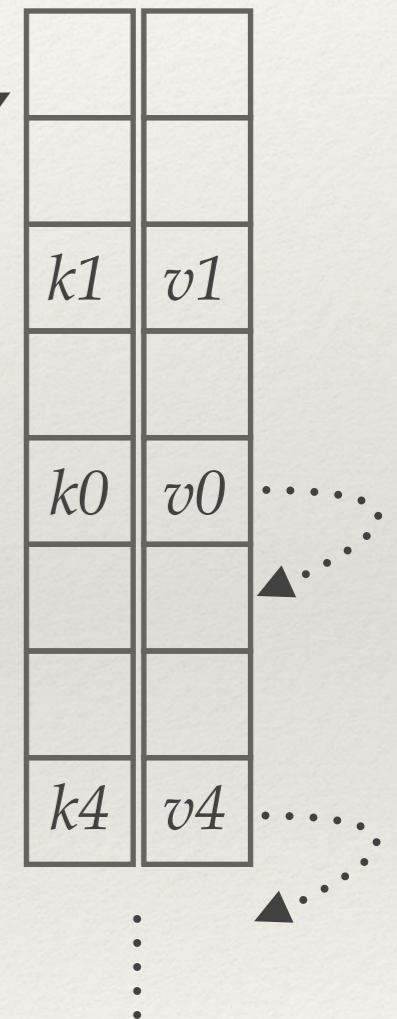
non-empty bucket **lanes**

# Hash Tables

❖ Vectorized hash tables

    ❖ *Vertical* vectorization

        ❖ W *input* keys at a time

        ❖ 1 *table* keys per input key

        ❖ *Gather* buckets

    ❖ Probing *(linear probing)*

        ❖ *Store* matching lanes

        ❖ *Replace* finished lanes

        ❖ *Keep* unfinished lanes

    ❖ Building *(linear probing)*

        ❖ *Keep* non-empty lanes

        ❖ *Detect* scatter *conflicts*

        ❖ *Scatter* empty lanes

linear probing hash table

input keys

$k1$
$k2$
$k3$
$k4$

hash indexes

$h1$
$h2$
$h3$
$h4$

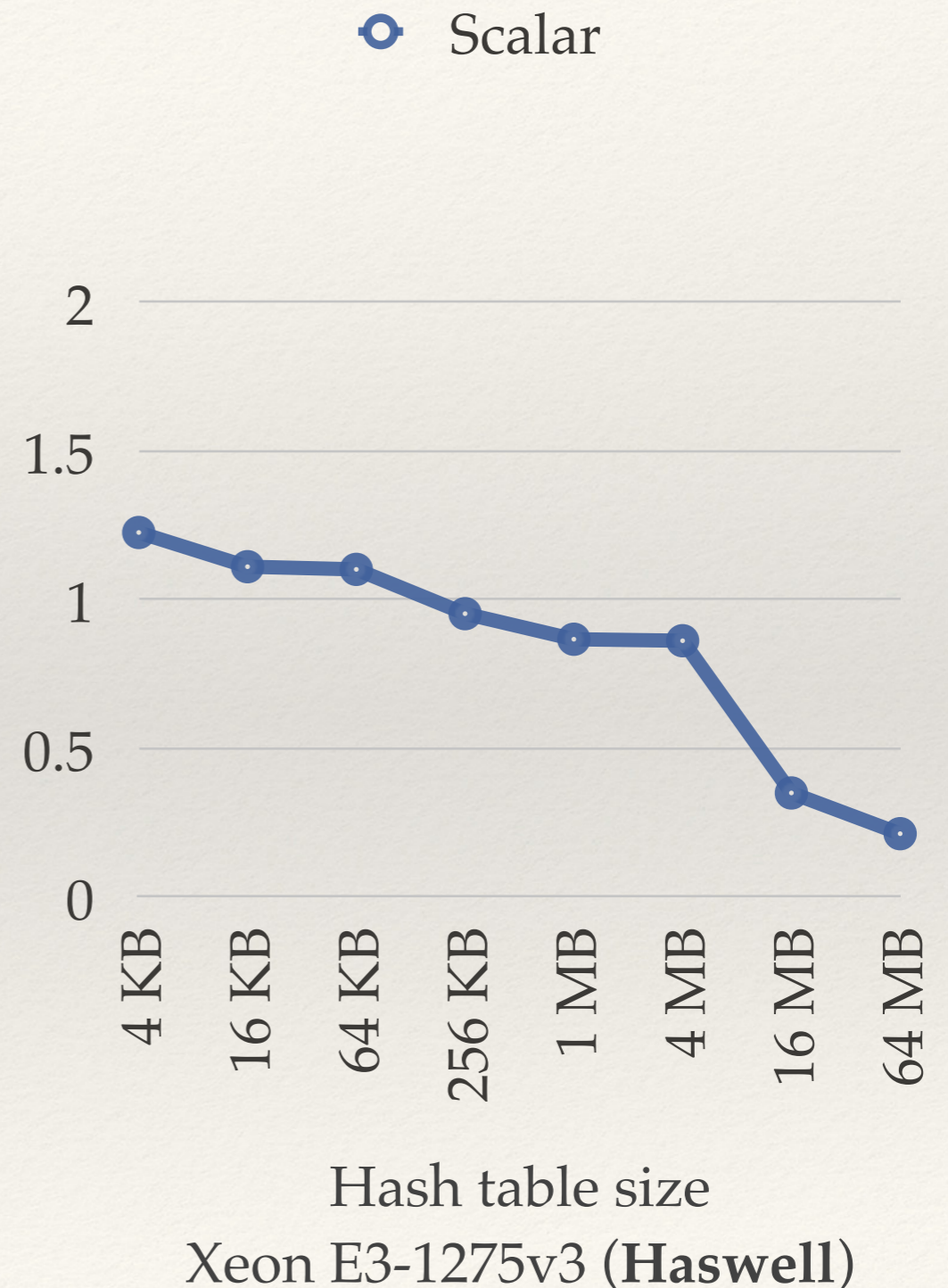*scatter tuples*

| $k1$ | $v1$ |
| $k0$ | $v0$ |
| $k4$ | $v4$ |

*conflicting lanes*

*non-conflicting lanes scattered*

*non-empty bucket lanes*

# Hash Tables

- Vectorized hash tables

  - *Vertical* vectorization

    - W *input* keys at a time

    - 1 *table* keys per input key

    - *Gather* buckets

  - Probing *(linear probing)*

    - *Store* matching lanes

    - *Replace* finished lanes

    - *Keep* unfinished lanes

  - Building *(linear probing)*

    - *Keep* non-empty lanes

    - *Detect* scatter *conflicts*

    - *Scatter* empty lanes & *replace*

    - *Keep* conflicting lanes

*selectively load*
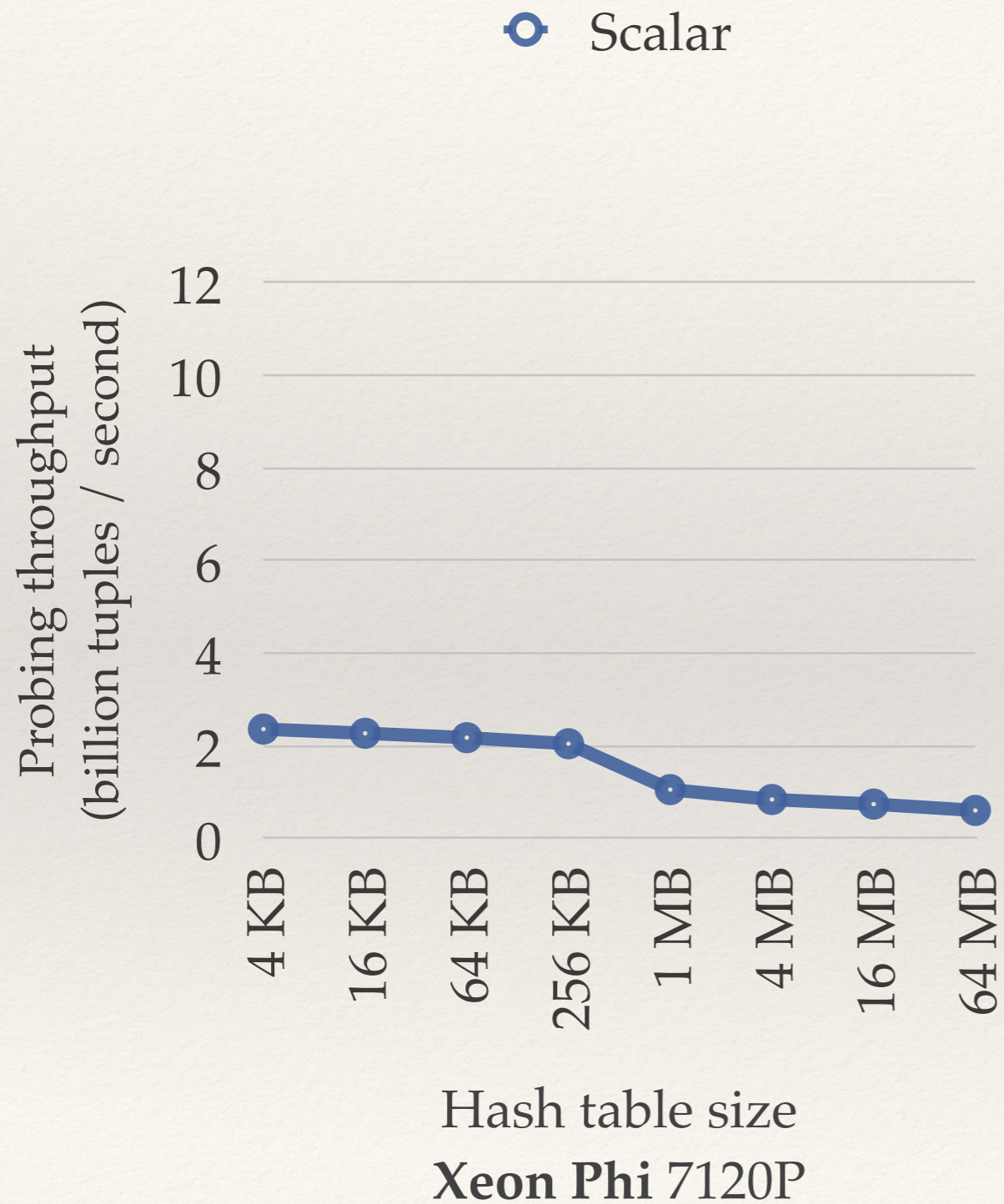*input keys*
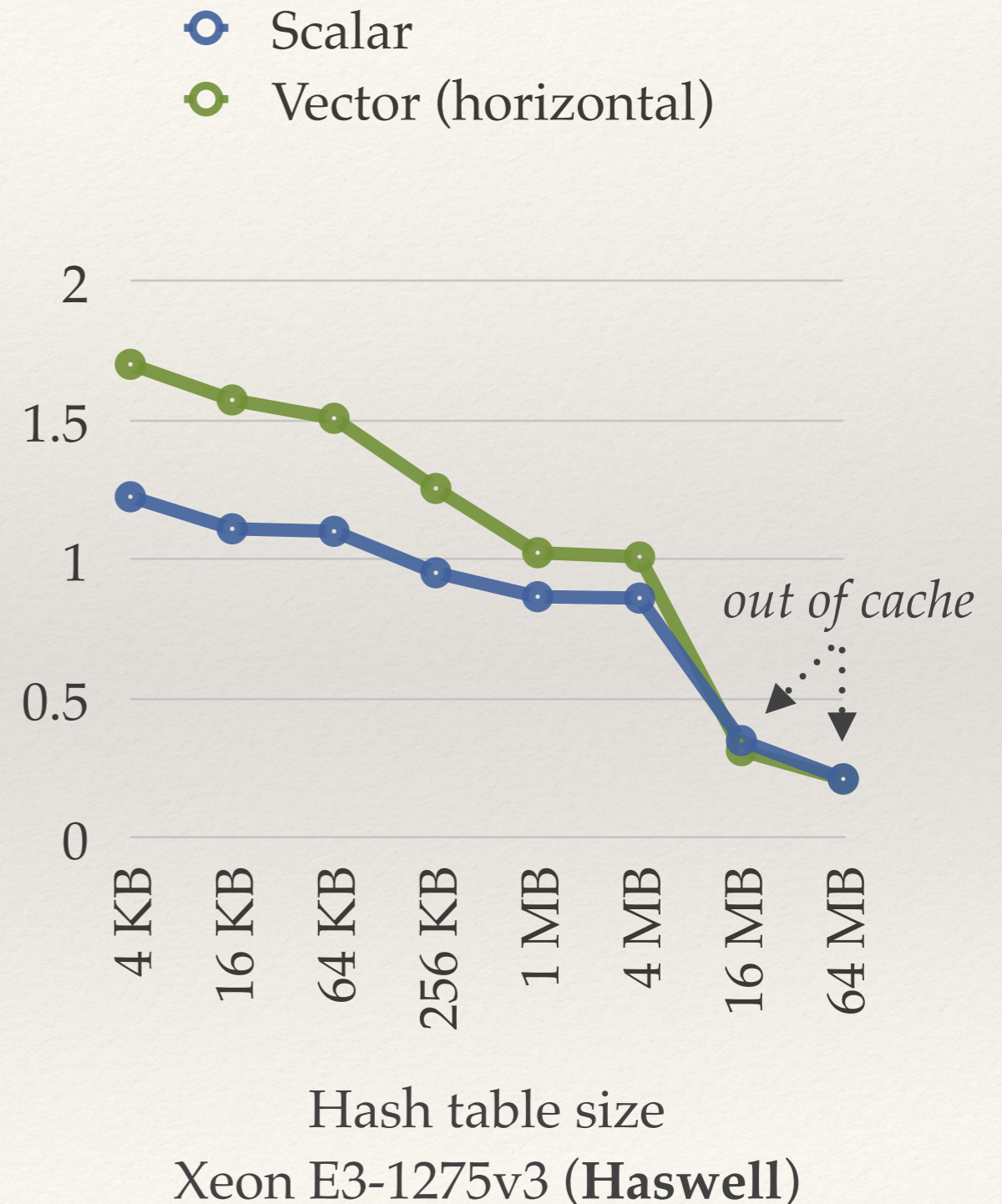*(2nd iteration)*

linear probing
hash table

| input keys | hash indexes |
|:---:|:---:|
| **k5** | h5 |
| k2 | h2 |
| k3 | h3 |
| **k6** | h6 |

| | |
|:---:|:---:|
| | |
| | |
| k1 | v1 |
| | |
| k0 | v0 |
| | |
| | |
| k4 | v4 |

*conflicting* **lanes** *kept*

*non-conflicting* **lanes** *replaced*

*non-empty bucket* **lanes** *kept*

# Hash Table (Linear Probing)



Scalar

Scalar

Probing throughput (billion tuples / second)

12
10
8
6
4
2
0

4 KB  16 KB  64 KB  256 KB  1 MB  4 MB  16 MB  64 MB

Hash table size
**Xeon Phi** 7120P

2
1.5
1
0.5
0

4 KB  16 KB  64 KB  256 KB  1 MB  4 MB  16 MB  64 MB

Hash table size
Xeon E3-1275v3 (**Haswell**)

# Hash Table (Linear Probing)

# Hash Table (Linear Probing)



Xeon Phi 7120P

Xeon E3-1275v3 (**Haswell**)

# Partitioning

- ❖ Types
  - ❖ Radix
    - ❖ 2 shifts  (in SIMD)
  - ❖ Hash
    - ❖ 2 multiplications  (in SIMD)
  - ❖ Range
    - ❖ Range function index
    - ❖ Binary search  (in the paper …)

# Partitioning

- ❖ Types
  - ❖ Radix
    - ❖ 2 shifts (in SIMD)
  - ❖ Hash
    - ❖ 2 multiplications (in SIMD)
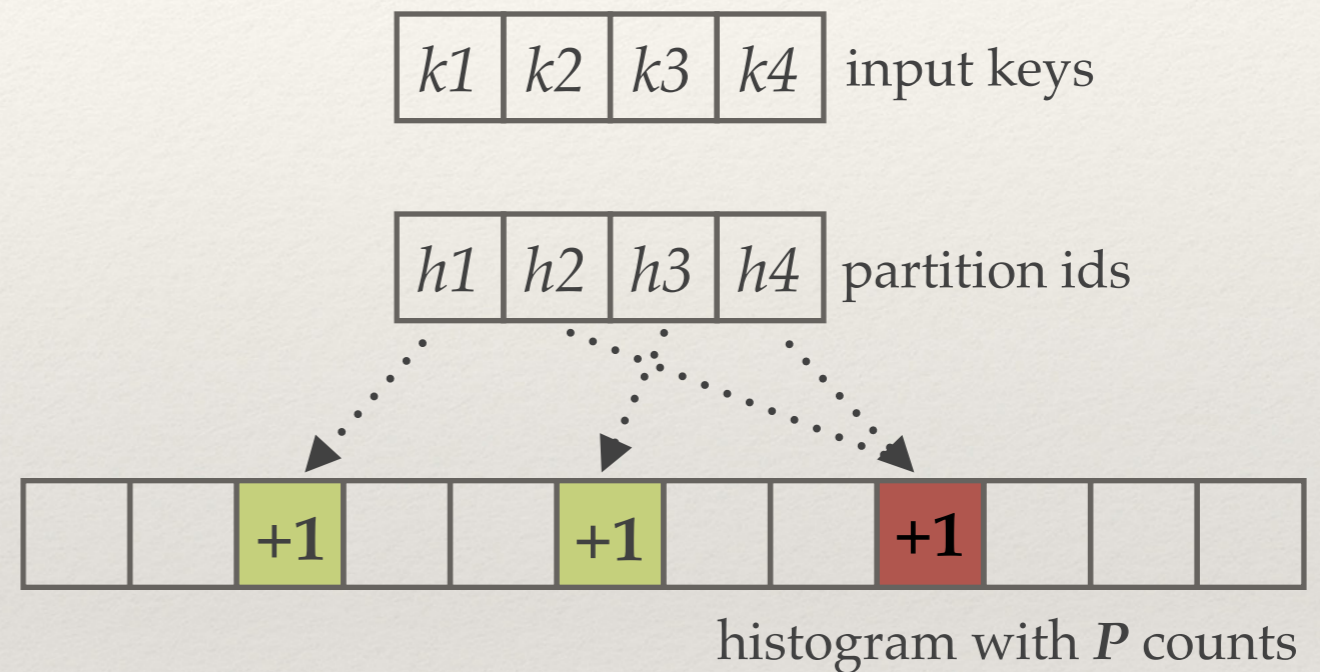  - ❖ Range
    - ❖ Range function index
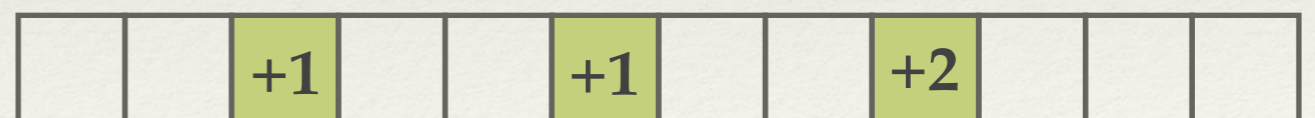    - ❖ Binary search (in the paper …)

- ❖ Histogram
  - ❖ Data parallel update
    - ❖ Gather & scatter counts
    - ❖ *Conflicts* miss counts

| $k1$ | $k2$ | $k3$ | $k4$ | input keys |

| $h1$ | $h2$ | $h3$ | $h4$ | partition ids |

| | | +1 | | | +1 | | | +1 | | | |

histogram with *P* counts

*… instead of …*

| | | +1 | | | +1 | | | +2 | | | |

# Partitioning

- ❖ Types
  - ❖ Radix
    - ❖ 2 shifts  (in SIMD)
  - ❖ Hash
    - ❖ 2 multiplications  (in SIMD)
  - ❖ Range
    - ❖ Range function index
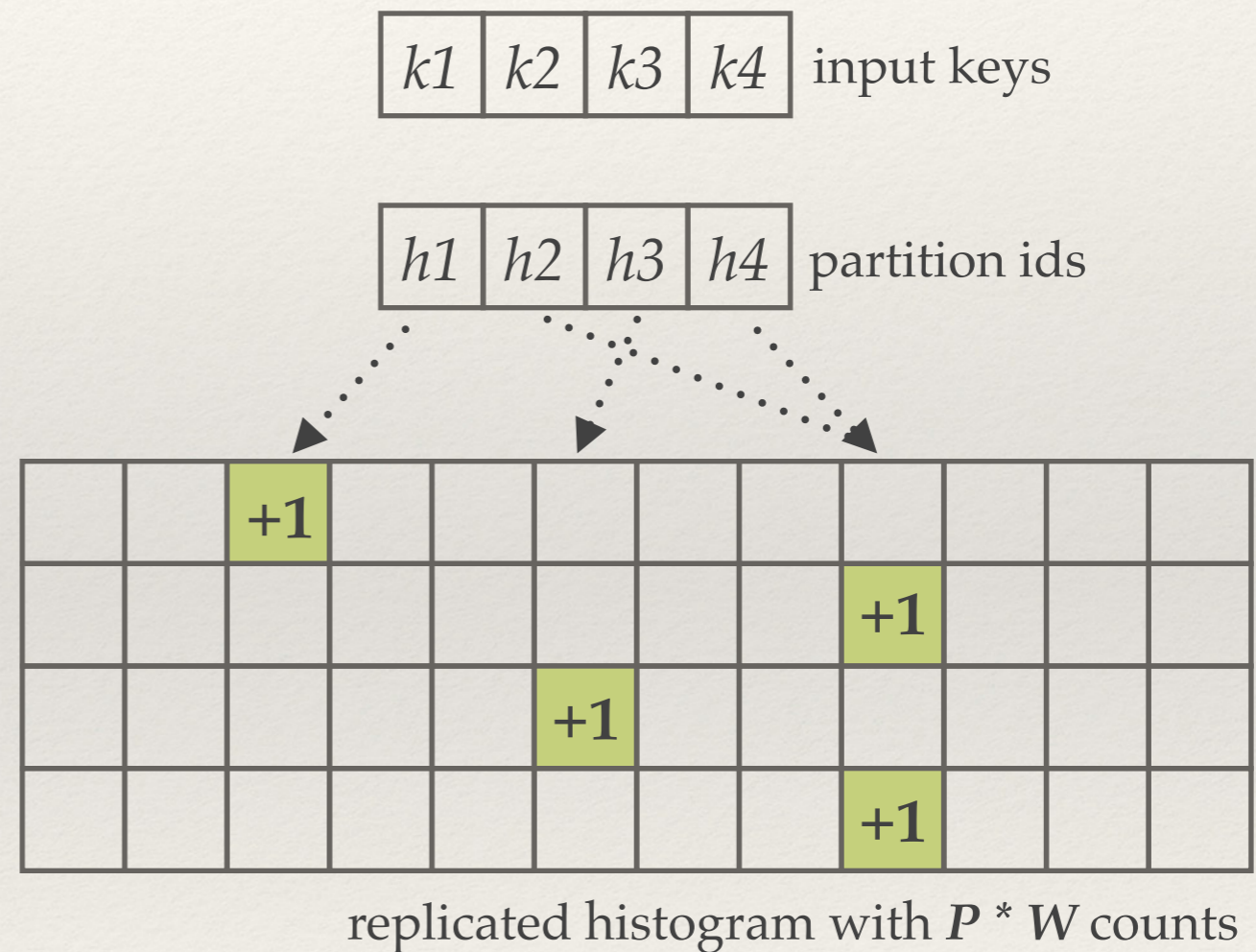    - ❖ Binary search  (in the paper …)

- ❖ Histogram
  - ❖ Data parallel update
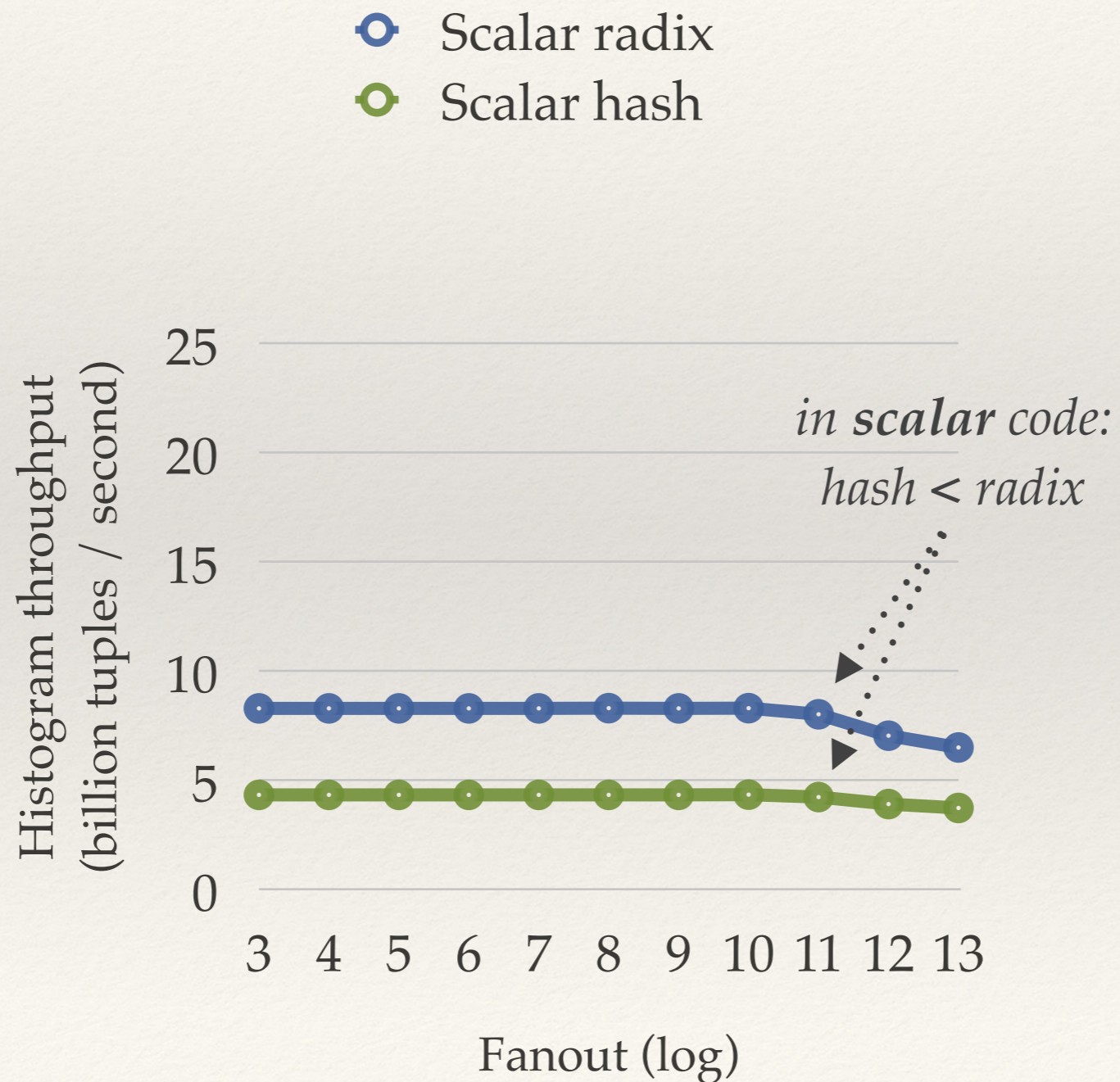    - ❖ Gather & scatter counts
    - ❖ *Conflicts* miss counts
    - ❖ *Replicate* histogram *W* times
    - ❖ More solutions in the paper …

| $k1$ | $k2$ | $k3$ | $k4$ | input keys |

| $h1$ | $h2$ | $h3$ | $h4$ | partition ids |

replicated histogram with *P * W* counts
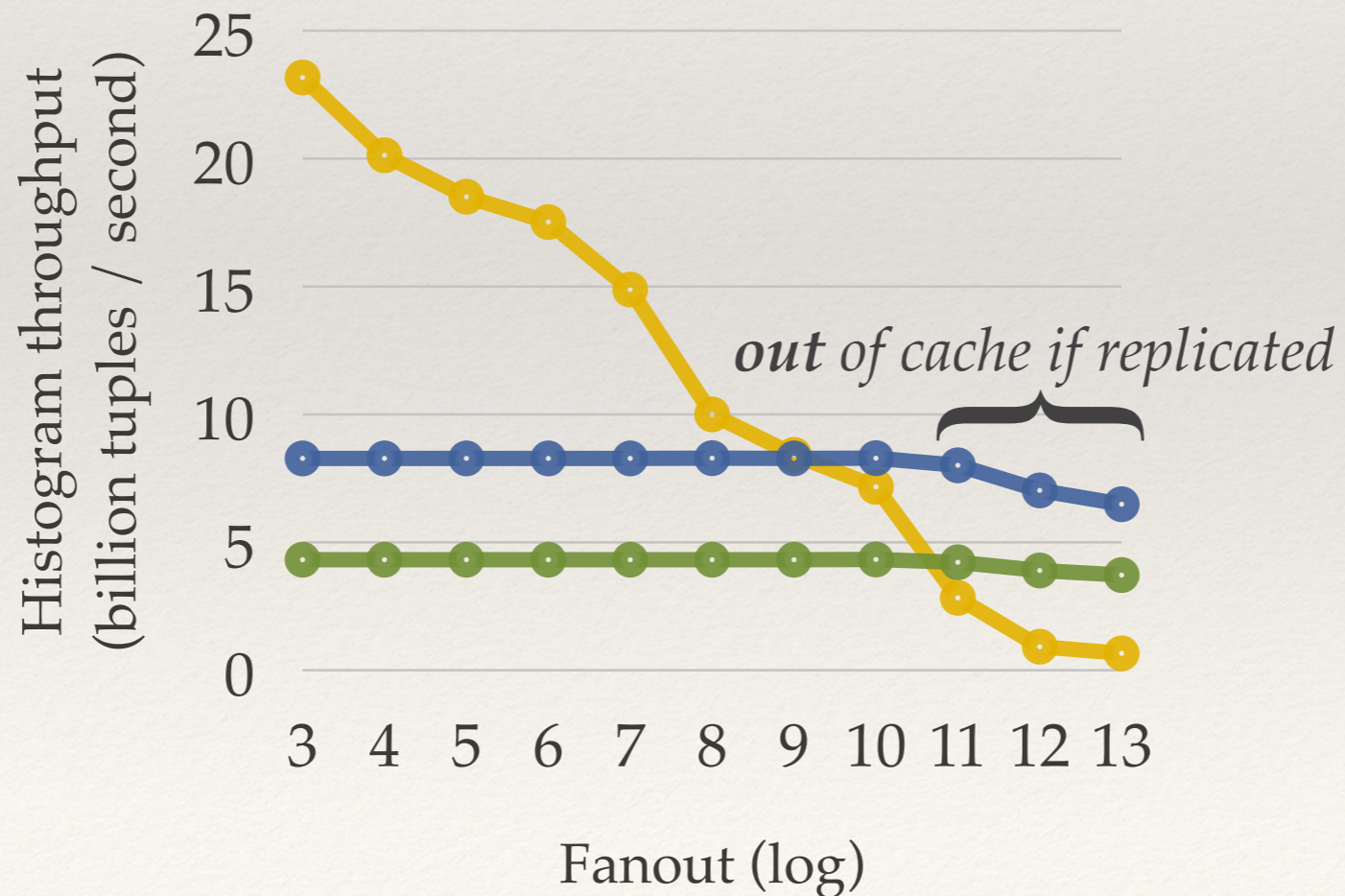
# Radix & Hash Histogram

Histogram generation on **Xeon Phi**
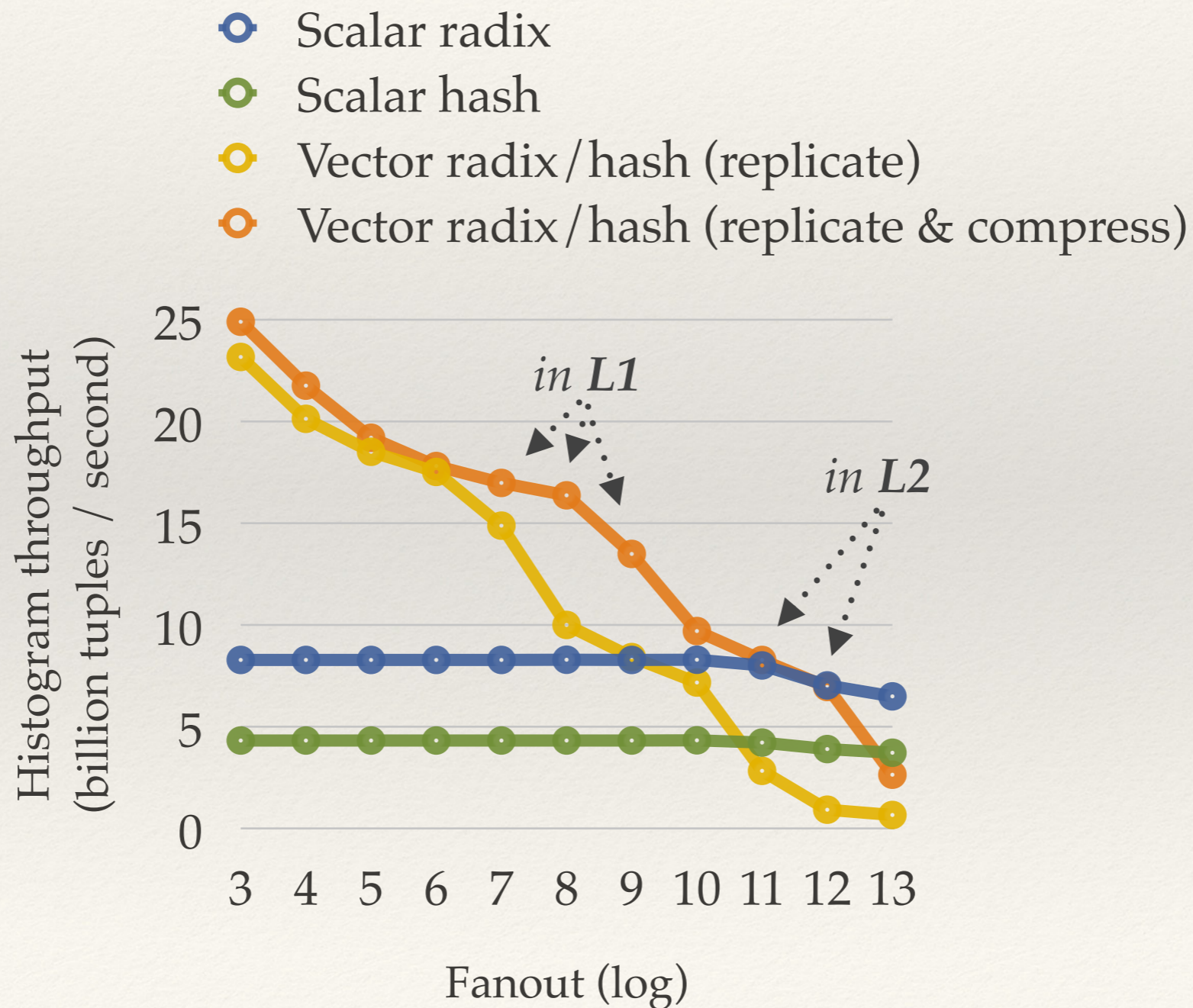
# Radix & Hash Histogram

Histogram generation on **Xeon Phi**

○ Scalar radix
○ Scalar hash
○ Vector radix / hash (replicate)
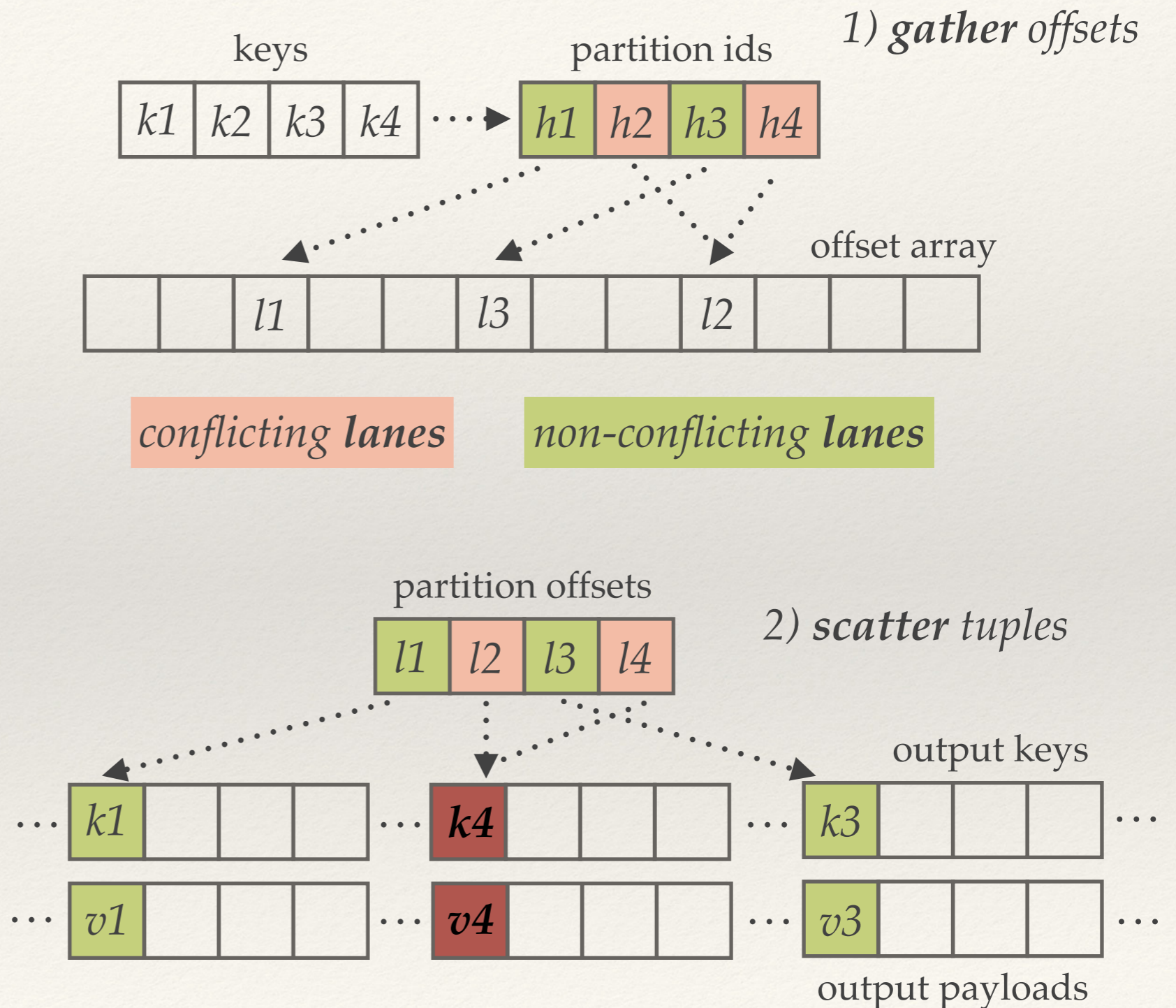
# Radix & Hash Histogram

Histogram generation on **Xeon Phi**

# Partitioning

❖ Shuffling

    ❖ Update the offsets

        ❖ Gather & scatter counts

        ❖ *Conflicts* miss counts

    ❖ Tuple transfer

        ❖ Scatter directly to *output*

        ❖ *Conflicts* overwrite tuples

# Partitioning

- Shuffling
  - Update the offsets
    - Gather & scatter counts
    - *Conflicts* miss counts
  - Tuple transfer
    - Scatter directly to output
    - *Conflicts* overwrite tuples
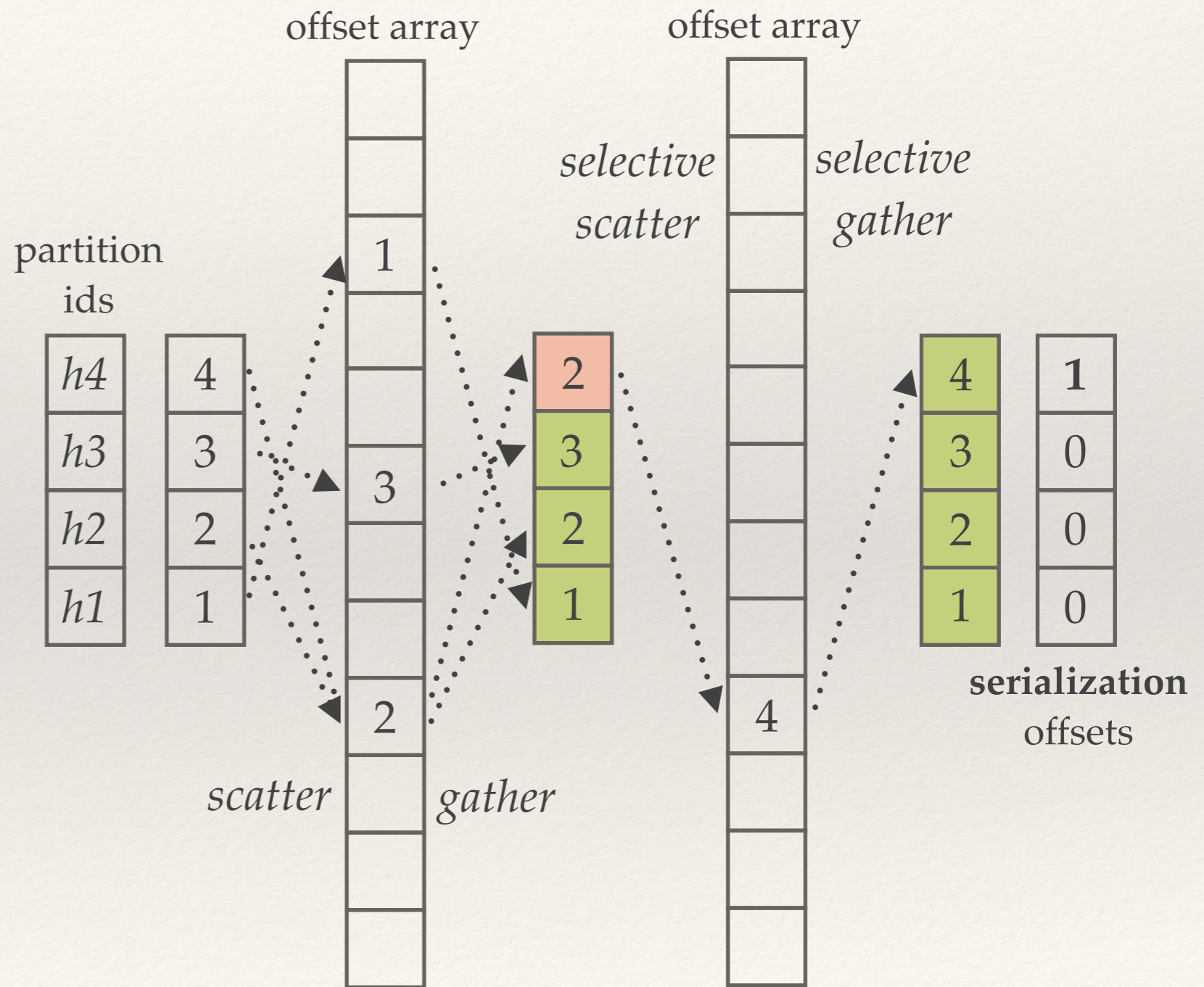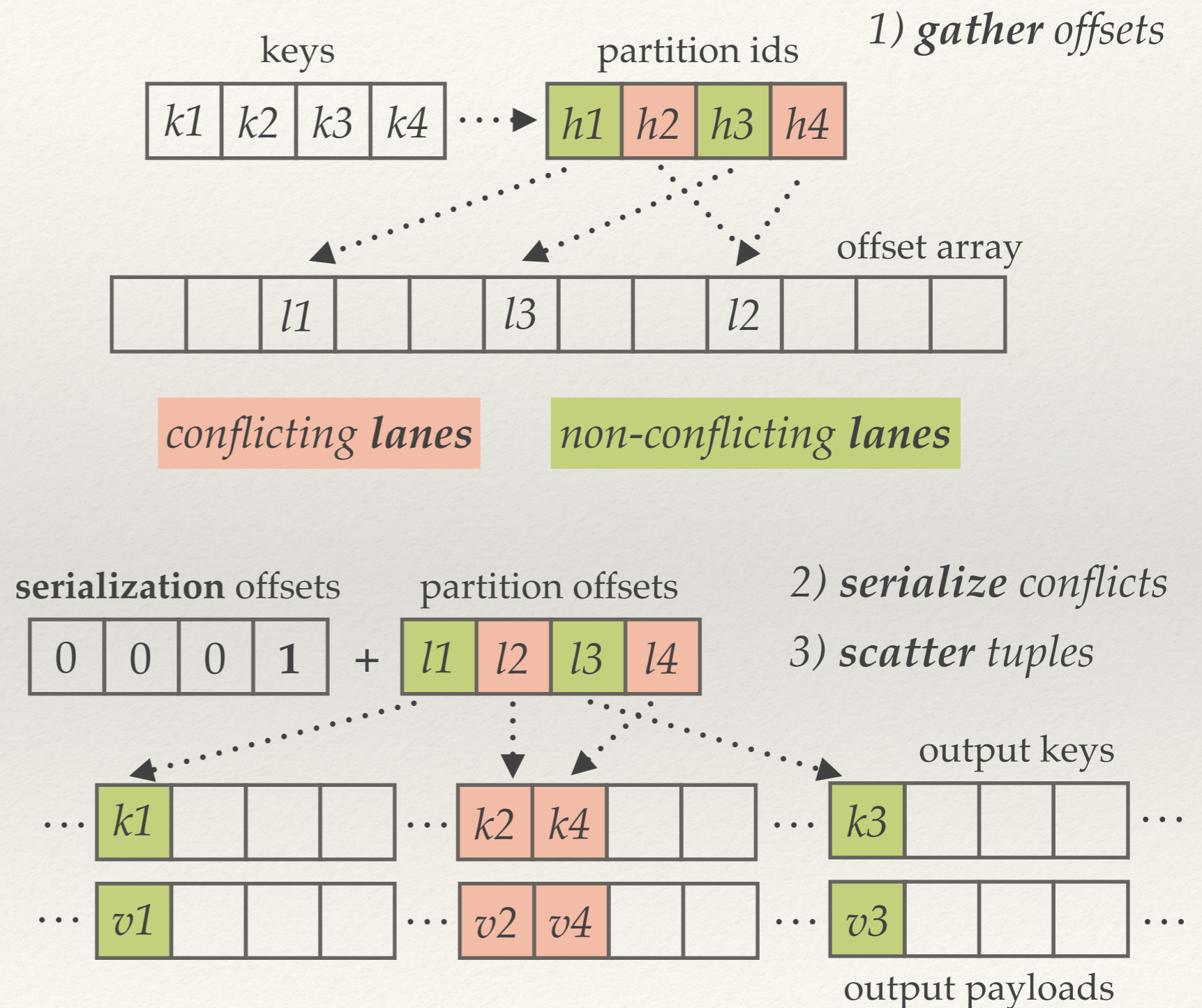    - *Serialize* conflicts

# Partitioning

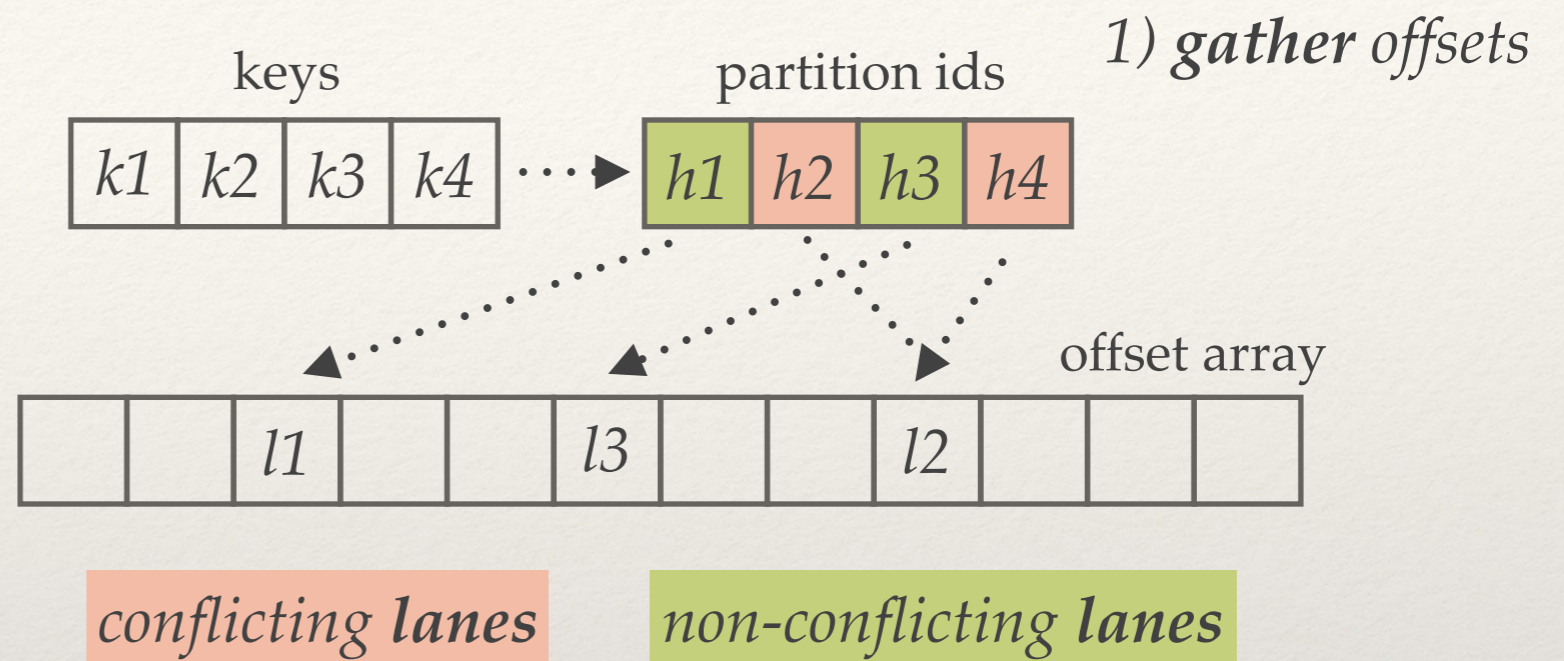- Shuffling
  - Update the offsets
    - Gather & scatter counts
    - *Conflicts* miss counts
  - Tuple transfer
    - Scatter directly to output
    - *Conflicts* overwrite tuples
    - *Serialize* conflicts
    - More in the paper …

*1) **gather** offsets*

keys

| k1 | k2 | k3 | k4 |
|----|----|----|----|

partition ids

| h1 | h2 | h3 | h4 |
|----|----|----|----|

offset array

|  |  | l1 |  | l3 |  | l2 |  |  |  |
|--|--|----|--|----|--|----|--|--|--|

*conflicting **lanes***     *non-conflicting **lanes***

**serialization** offsets

| 0 | 0 | 0 | **1** |
|---|---|---|---|

**+**

partition offsets

| l1 | l2 | l3 | l4 |
|----|----|----|----|

*2) **serialize** conflicts*

*3) **scatter** tuples*

output keys

… | k1 |  |  |  | …   … | k2 | k4 |  |  | …   … | k3 |  |  |  | …

… | v1 |  |  |  | …   … | v2 | v4 |  |  | …   … | v3 |  |  |  | …

output payloads

# Partitioning
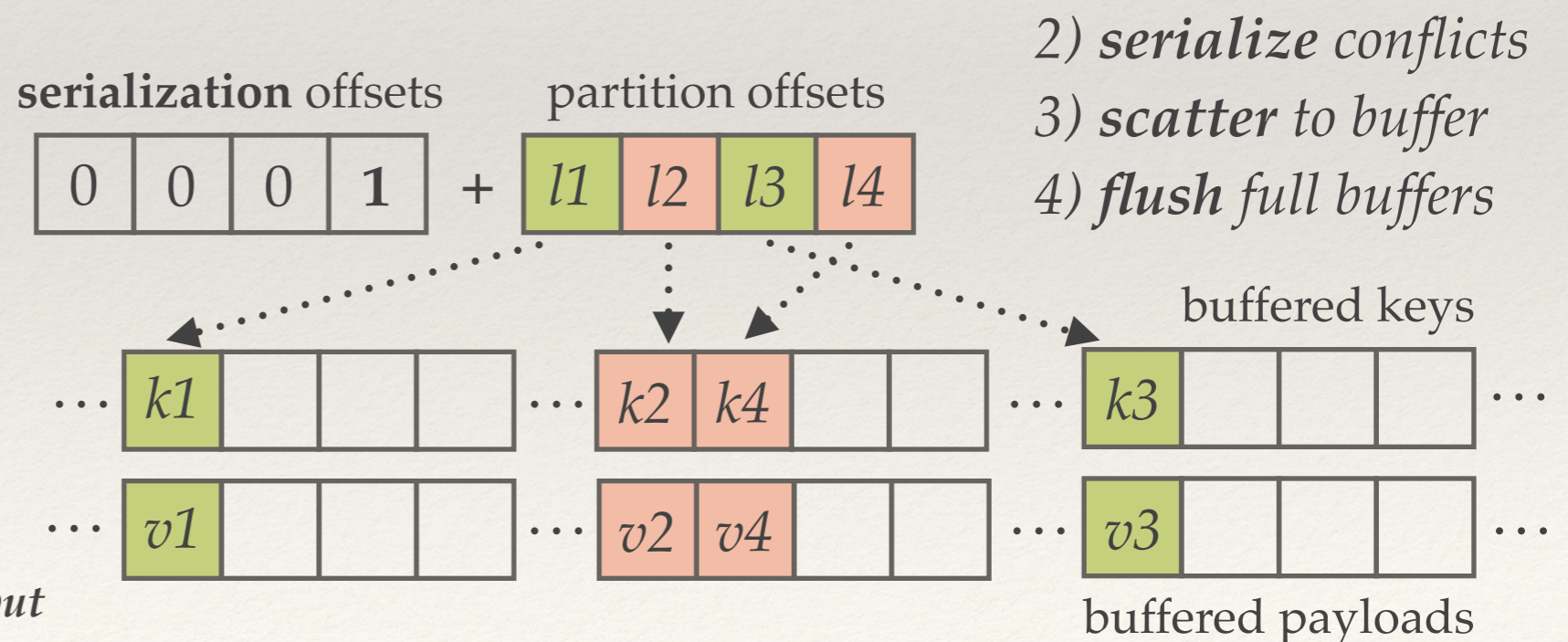
- Shuffling
  - Update the offsets
    - Gather & scatter counts
    - *Conflicts* miss counts
  - Tuple transfer
    - Scatter directly to output
    - *Conflicts* overwrite tuples
    - *Serialize* conflicts
    - More in the paper …
  - Buffering
    - When input >> cache
    - Handle *conflicts*
    - Scatter tuples to *buffers*
    - Buffers are *cache*-resident
    - Stream *full* buffers to *output*

*1) **gather** offsets*

keys

| k1 | k2 | k3 | k4 |

partition ids

| h1 | h2 | h3 | h4 |

offset array

| | | l1 | | l3 | | l2 | | | |

*conflicting **lanes**    non-conflicting **lanes***

*2) **serialize** conflicts*
*3) **scatter** to buffer*
*4) **flush** full buffers*

**serialization** offsets

| 0 | 0 | 0 | **1** |

+

partition offsets

| l1 | l2 | l3 | l4 |

buffered keys

… | k1 | | | | …    … | k2 | k4 | | | …    … | k3 | | | | …

… | v1 | | | | …    … | v2 | v4 | | | …    … | v3 | | | | …

buffered payloads

# Partitioning

- Shuffling
  - Update the offsets
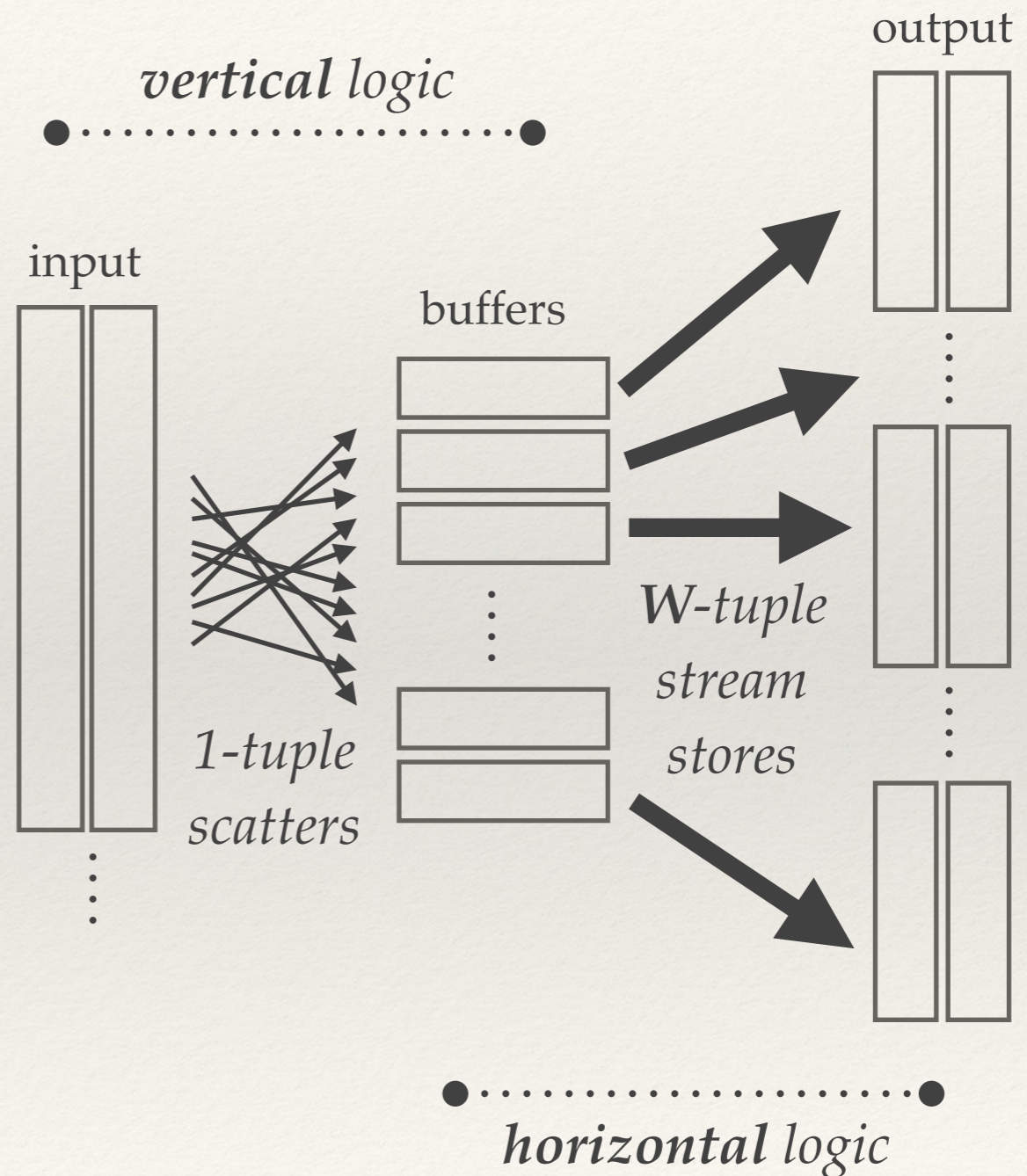    - Gather & scatter counts
    - *Conflicts* miss counts
  - Tuple transfer
    - Scatter directly to output
    - *Conflicts* overwrite tuples
    - *Serialize* conflicts
    - More in the paper …
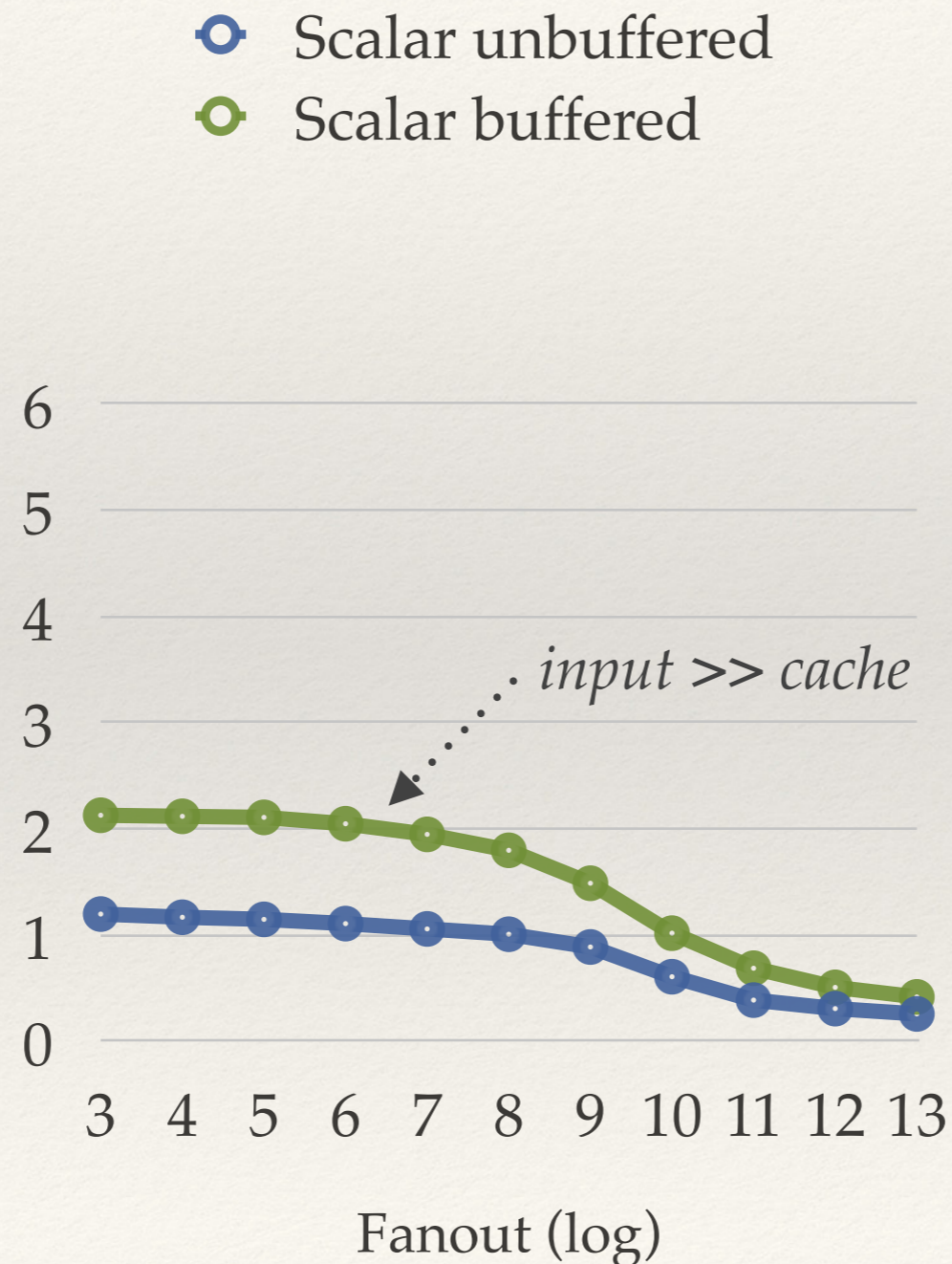  - Buffering
    - When input >> cache
    - Handle *conflicts*
    - Scatter tuples to *buffers*
    - Buffers are *cache*-resident
    - Stream *full* buffers to *output*
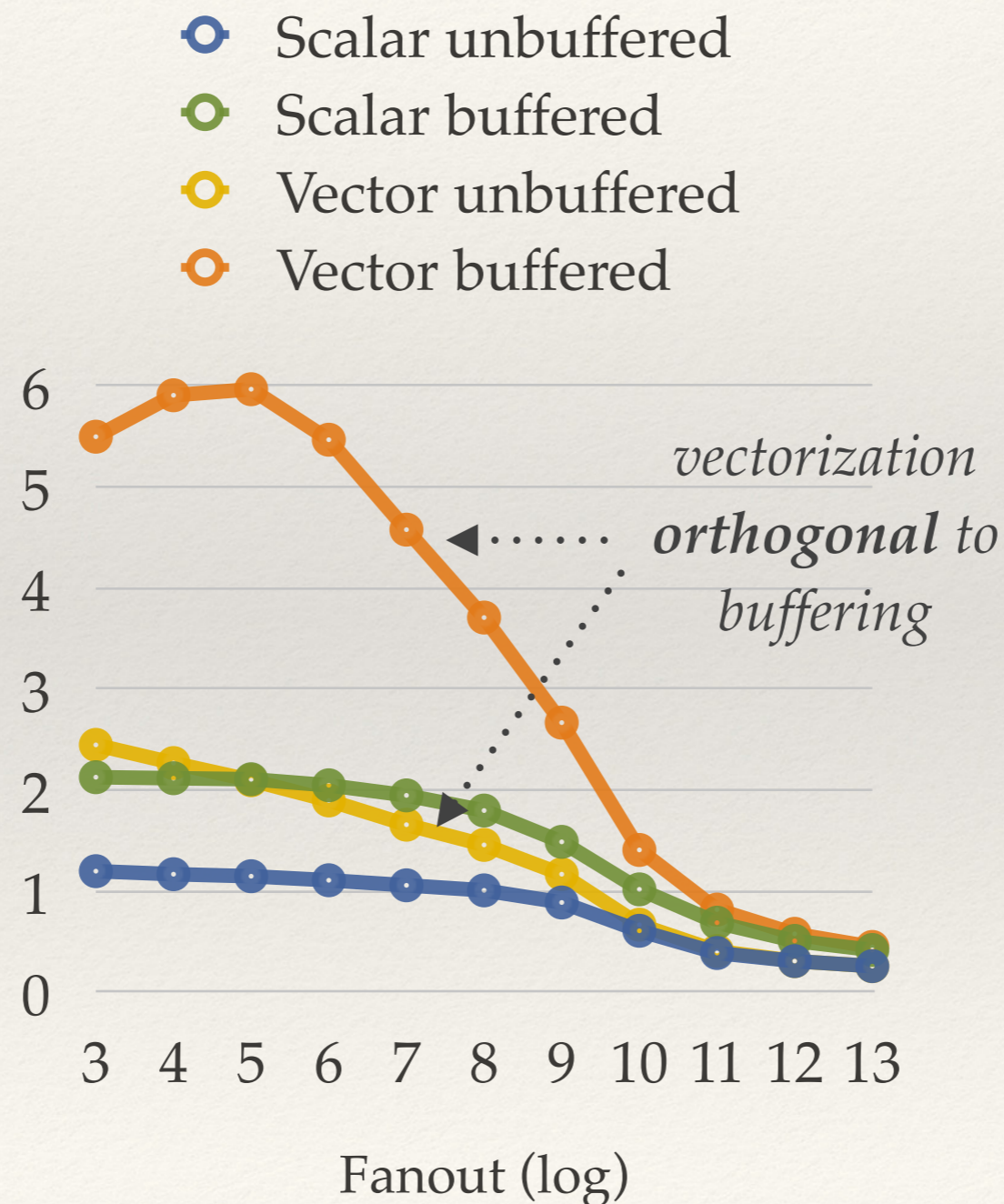
*vertical* logic

input

buffers

output

*1-tuple scatters*

*W-tuple stream stores*

*horizontal* logic

# Buffered & Unbuffered Partitioning

Large-scale data shuffling on **Xeon Phi**

○ Scalar unbuffered
○ Scalar buffered



*input >> cache*

Fanout (log)

# Buffered & Unbuffered Partitioning

Large-scale data shuffling on **Xeon Phi**

- ○ Scalar unbuffered
- ○ Scalar buffered
- ○ Vector unbuffered
- ○ Vector buffered

*vectorization* **orthogonal** *to buffering*

Fanout (log)

# Sorting & Hash Join Algorithms

❖ Least-significant-bit (LSB) radix-sort

    ❖ *Stable* radix partitioning passes

        ❖ *Fully* vectorized

# Sorting & Hash Join Algorithms

❖ Least-significant-bit (LSB) radix-sort

    ❖ *Stable* radix partitioning passes

        ❖ *Fully* vectorized

❖ Hash join

    ❖ *No* partitioning

        ❖ Build 1 *shared* hash table using *atomics*

        ❖ *Partially* vectorized

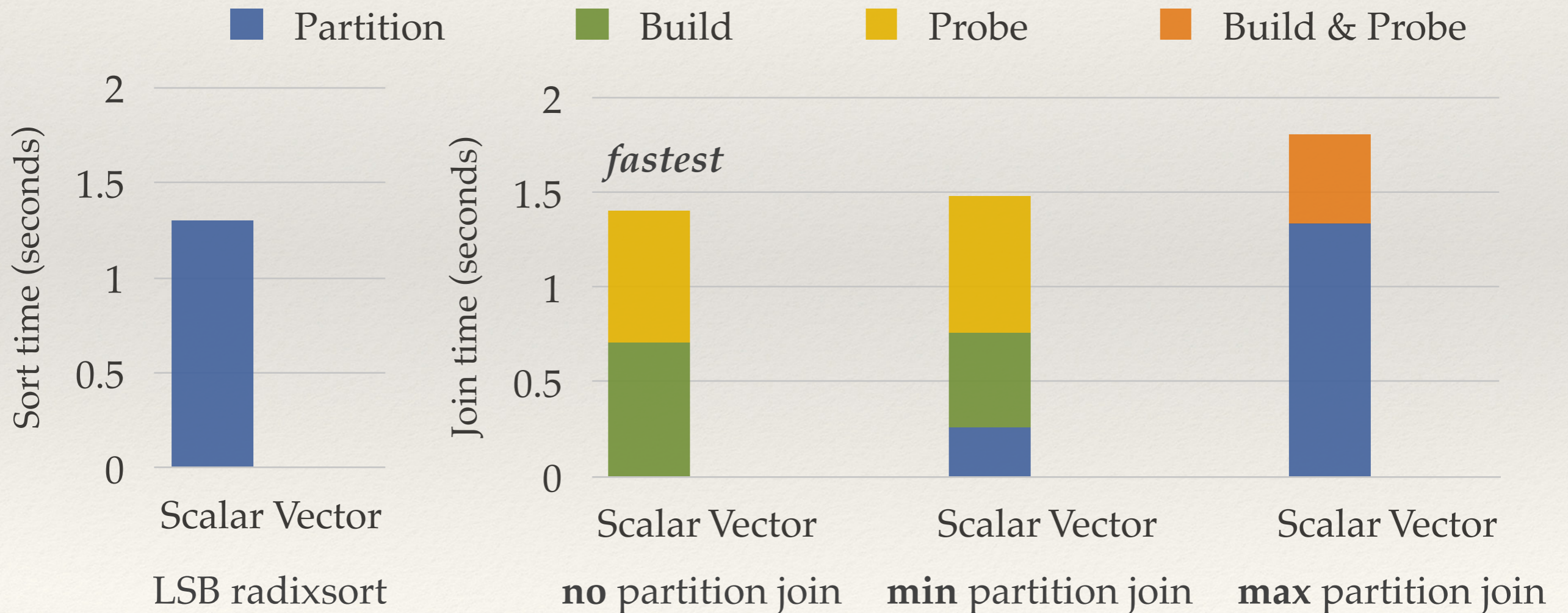# Sorting & Hash Join Algorithms

❖ Least-significant-bit (LSB) radix-sort

    ❖ *Stable* radix partitioning passes

        ❖ *Fully* vectorized

❖ Hash join

    ❖ *No* partitioning

        ❖ Build 1 *shared* hash table using *atomics*

        ❖ *Partially* vectorized

    ❖ *Min* partitioning

        ❖ Partition *building* table

        ❖ Build 1 hash table per thread

        ❖ *Fully* vectorized

# Sorting & Hash Join Algorithms

* Least-significant-bit (LSB) radix-sort
    * *Stable* radix partitioning passes
        * *Fully* vectorized

* Hash join
    * *No* partitioning
        * Build 1 *shared* hash table using *atomics*
        * *Partially* vectorized
    * *Min* partitioning
        * Partition *building* table
        * Build 1 hash table per thread
        * *Fully* vectorized
    * *Max* partitioning
        * Partition *both* tables repeatedly
        * Build & probe *cache-resident* hash tables
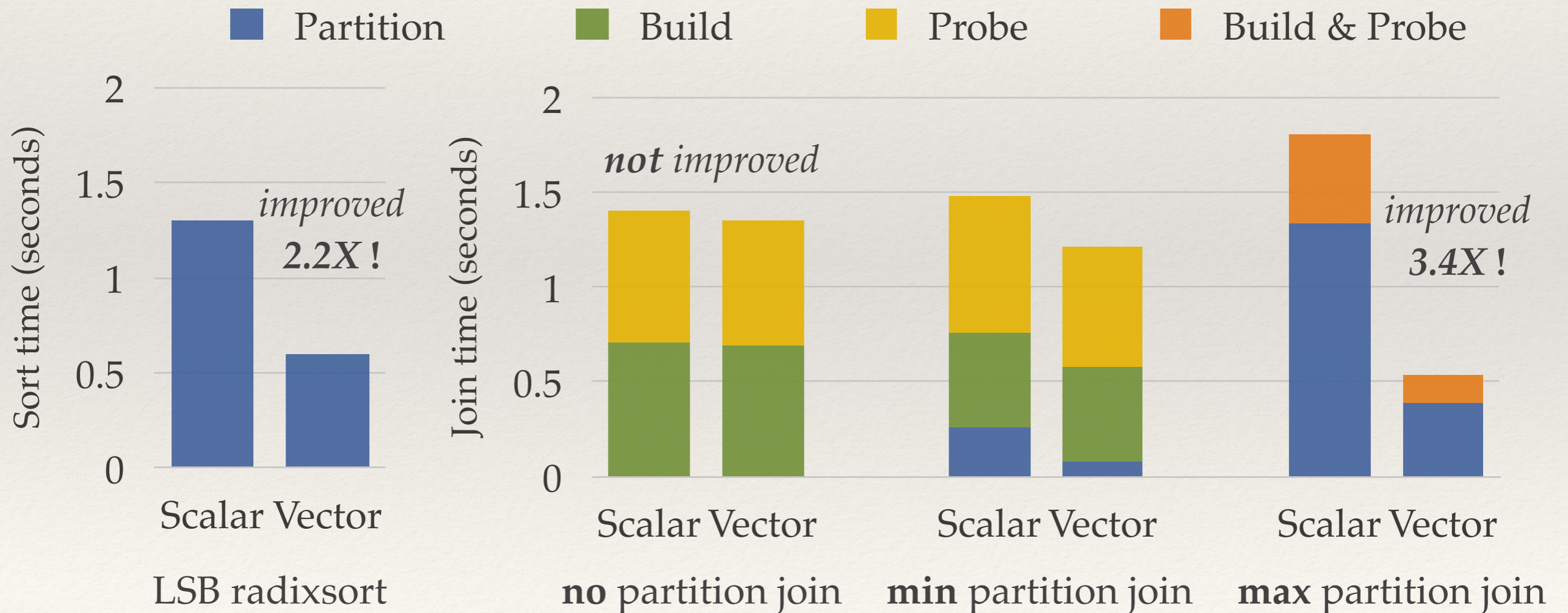        * *Fully* vectorized

# Hash Joins

**sort** 400 million tuples &
**join** 200 & 200 million tuples
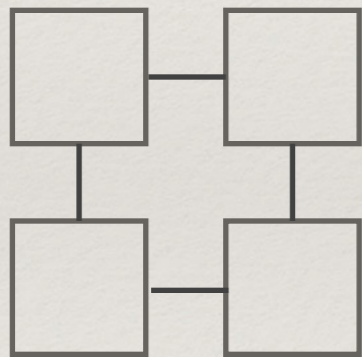32-bit keys & payloads
on **Xeon Phi**



Legend: ■ Partition  ■ Build  ■ Probe  ■ Build & Probe

Left chart: Sort time (seconds), Scalar Vector, LSB radixsort

Right chart: Join time (seconds), *fastest*, Scalar Vector — **no** partition join, **min** partition join, **max** partition join

# Hash Joins

sort 400 million tuples &
join 200 & 200 million tuples
32-bit keys & payloads
on Xeon Phi

# Power Efficiency

4 CPUs
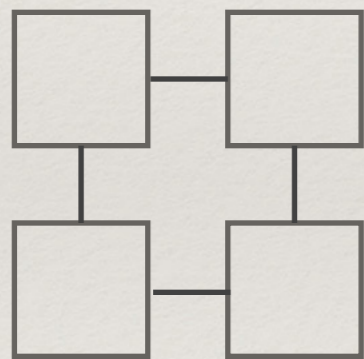(Xeon E5-4620)
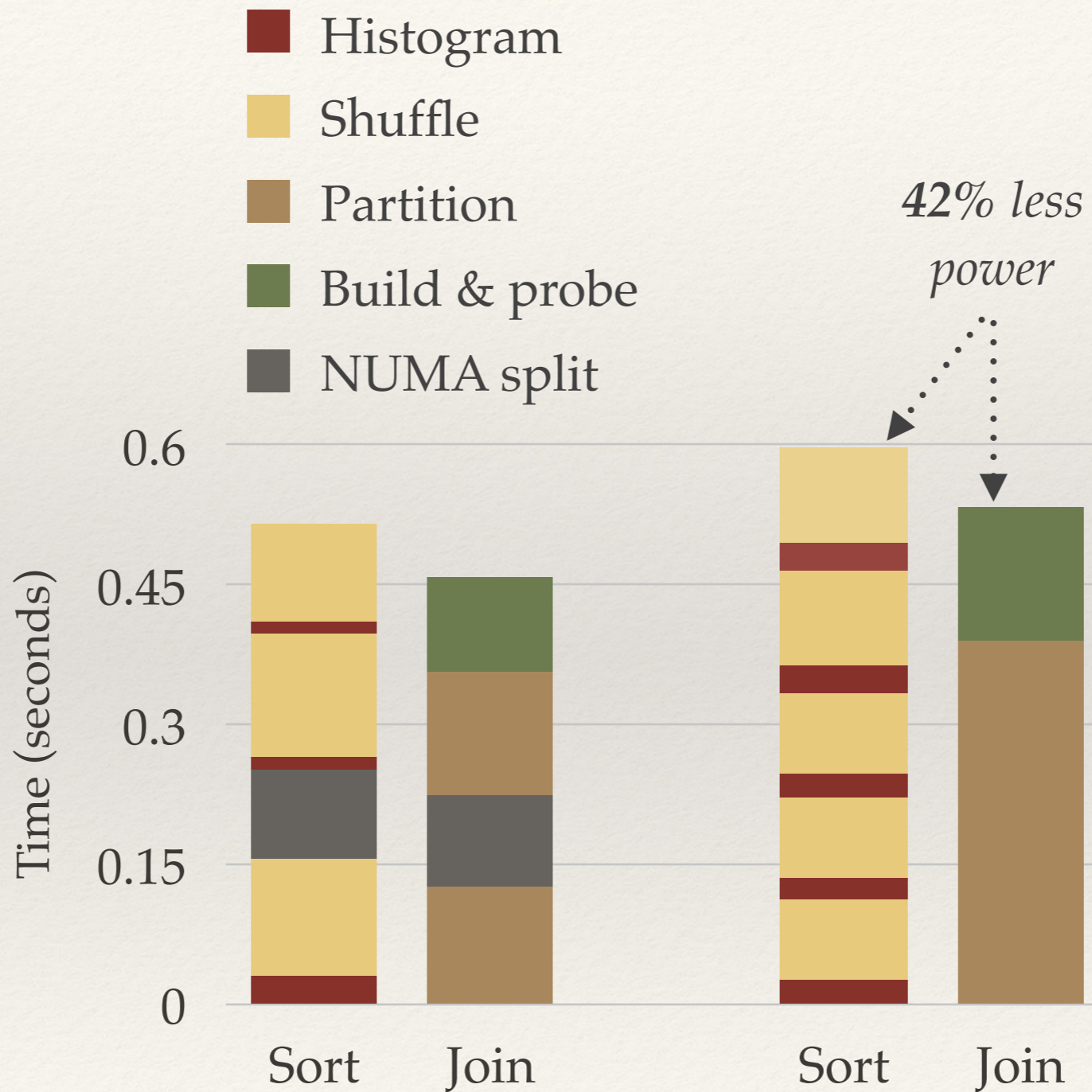32 **Sandy**
**Bridge** cores
**300** Watts TDP

**2.2** GHz
**40** GB/s

Histogram
Shuffle
Partition
Build & probe
NUMA split

# Conclusions

❖ *Full* vectorization

   ❖ $O(f(n))$ scalar —> $O(f(n)/W)$ vector operations

   ❖ Good vectorization principles improve *performance*

      ❖ Define & *reuse* fundamental operations

      ❖ e.g. *vertical* vectorization, maximize lane *utilization*, …

# Conclusions

❖ *Full* vectorization

   ❖ $O(f(n))$ scalar —> $O(f(n)/W)$ vector operations

   ❖ Good vectorization principles improve *performance*

      ❖ Define & *reuse* fundamental operations

      ❖ e.g. *vertical* vectorization, maximize lane *utilization*, …


❖ Impact on software design

   ❖ Vectorization favors cache-conscious algorithms

      ❖ e.g. *partitioned* hash join >> non-partitioned hash join if vectorized

   ❖ Vectorization is *orthogonal* to other optimizations

      ❖ e.g. *both* unbuffered & *buffered* partitioning get vectorization speedup
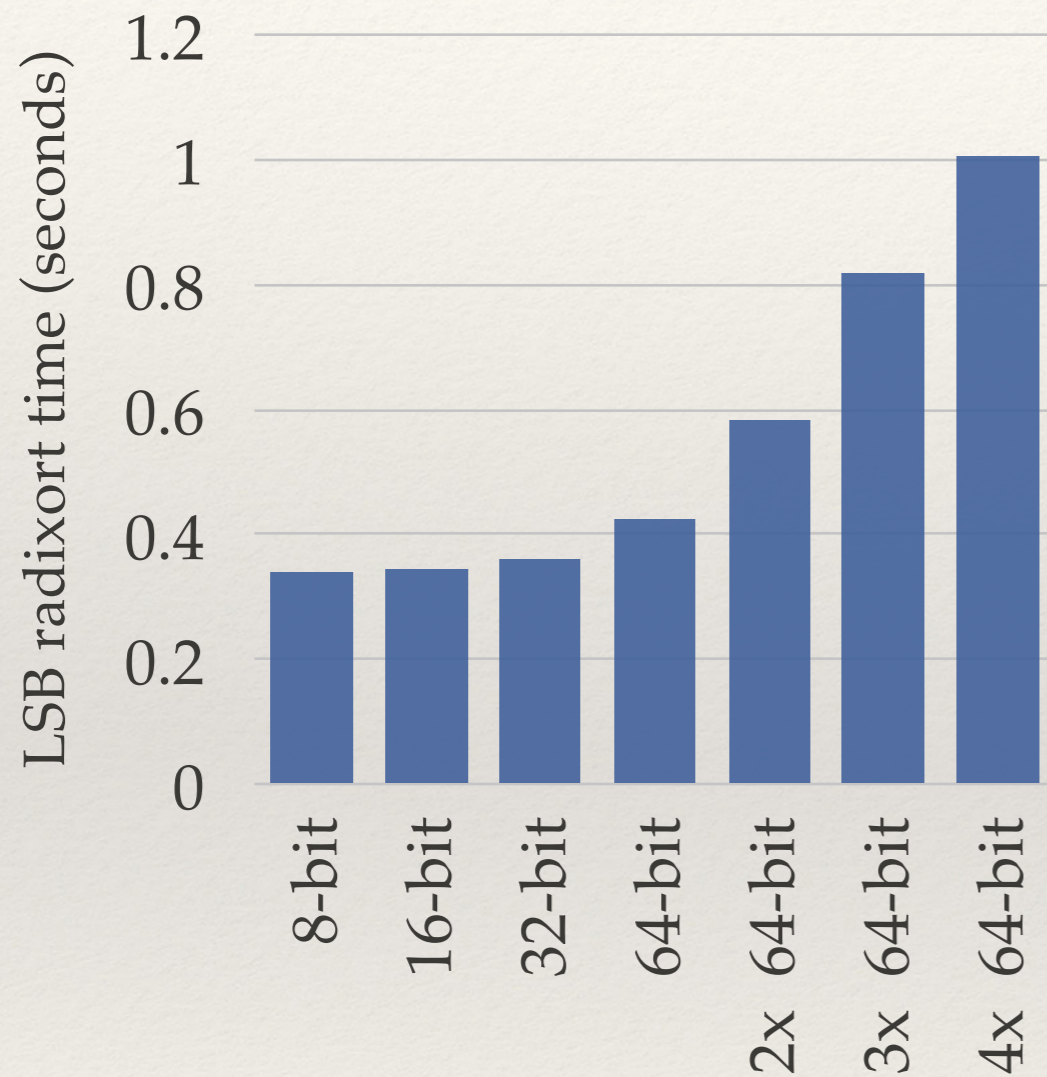
# Conclusions

- *Full* vectorization
  - $O(f(n))$ scalar —> $O(f(n)/W)$ vector operations
  - Good vectorization principles improve *performance*
    - Define & *reuse* fundamental operations
    - e.g. *vertical* vectorization, maximize lane *utilization*, …

- Impact on software design
  - Vectorization favors cache-conscious algorithms
    - e.g. *partitioned* hash join >> non-partitioned hash join if vectorized
  - Vectorization is *orthogonal* to other optimizations
    - e.g. *both* unbuffered & *buffered* partitioning get vectorization speedup

- Impact on hardware design
  - *Simple* cores almost as fast as *complex* cores  (for OLAP)
    - 61 simple *P54C* cores ~ 32 complex *Sandy Bridge* cores
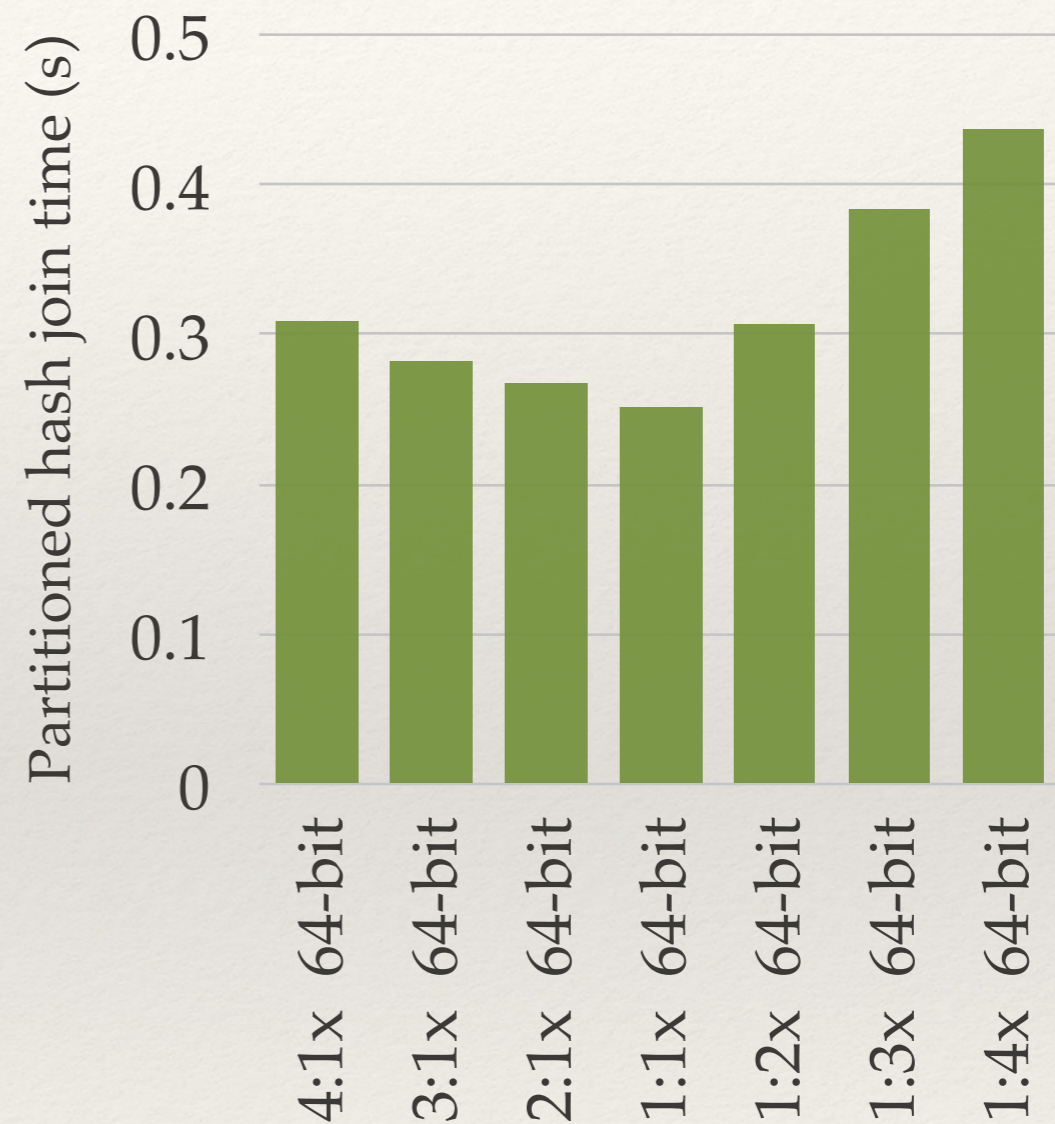  - Improved *power* efficiency for analytical databases

# Questions

# Join & Sort with Early Materialized Payloads



LSB radixort time (seconds) — 200 million 32-bit keys & X payloads on **Xeon Phi**

Bars: 8-bit, 16-bit, 32-bit, 64-bit, 2x 64-bit, 3x 64-bit, 4x 64-bit

Partitioned hash join time (s) — 10 & 100 million 32-bit keys & X payloads on **Xeon Phi**

Bars: 4:1x 64-bit, 3:1x 64-bit, 2:1x 64-bit, 1:1x 64-bit, 1:2x 64-bit, 1:3x 64-bit, 1:4x 64-bit